

1 **Within and among farm variability of coffee quality of smallholders in** 2 **southwest Ethiopia**

3
4 Merkebu Getachew^{ab}, Pascal Boeckx^c, Kris Verheyen^a, Kassaye Tolassa^d, Ayco J.M. Tack^e, Kristoffer
5 Hylander^e, Stijn Luca^f, Beyene Zewdie^e, Pieter De Frenne^a

6
7 ^aForest & Nature Lab, Department of Environment, Faculty of Bioscience Engineering, Ghent University,
8 Belgium

9 ^bDepartment of Horticulture and Plant Sciences, College of Agriculture and Veterinary Medicine, Jimma
10 University, Ethiopia

11 ^cIsotope Bioscience Laboratory, Department of Green Chemistry and Technology, Faculty of Bioscience
12 Engineering, Ghent University, Gent, Belgium

13 ^dEthiopian Institute of Agricultural Research, Addis Ababa, Ethiopia

14 ^eDepartment of Ecology, Environment and Plant Sciences, Stockholm University, Stockholm, Sweden

15 ^fDepartment of Data Analysis and Mathematical Modelling, Faculty of Bioscience Engineering, Ghent
16 University, Belgium.

17 18 19 **Abstract**

20 The biophysical drivers that affect coffee quality vary within and among farms. Quantifying their
21 relative importance is crucial for making informed decisions concerning farm management,
22 marketability and profit for coffee farmers. The present study was designed to quantify the relative
23 importance of biophysical variables affecting coffee bean quality within and among coffee farms
24 and to evaluate a near infrared spectroscopy-based model to predict coffee quality. Twelve coffee
25 plants growing under low, intermediate and dense shade were studied in twelve coffee farms across
26 an elevational gradient (1470–2325 m asl) in Ethiopia. We found large within farm variability,
27 demonstrating that conditions varying at the coffee plant-level are of large importance for physical
28 attributes and cupping scores of green coffee beans. Overall, elevation appeared to be the key
29 biophysical variable influencing all the measured coffee bean quality attributes at the farm level
30 while canopy cover appeared to be the most important biophysical variable driving the above-
31 mentioned coffee bean quality attributes at the coffee plant level. The biophysical variables driving
32 coffee quality (total preliminary and specialty quality) were the same as those driving variations
33 in the near-infrared spectroscopy data, which supports future use of this technology to assess green
34 bean coffee quality. Most importantly, our findings show that random forest is computationally
35 fast and robust to noise, besides having comparable prediction accuracy. Hence, it is a useful
36 machine learning tool for regression studies and has potential for modeling linear and nonlinear

37 multivariate calibrations. The study also confirmed that near-infrared spectroscopic-based
38 predictions can be applied as a supplementary approach for coffee cup quality evaluations.

39

40 **Keywords:** coffee quality, near-infrared spectroscopy (NIRS), random forest, PERMANOVA

41

42 **1. Introduction**

43

44 Coffee is one of the most important global commodities providing livelihood opportunities to
45 millions of people in the global South (Legesse, 2020; Davis *et al.* 2019; Ovalle-Rivera *et al.*
46 2015). In addition to being an important cash crop to farmers in Ethiopia, it generates about a
47 quarter of the country's export earnings (Legesse, 2020). Quality is becoming paramount in the
48 global coffee market. Coffee quality is about having desirable characteristics such as clean in its
49 appearance and good cupping scores (Carvalho *et al.* 2020). High-quality coffee shows little or no
50 physical defects, for instance, broken beans, insect damage, and other foreign materials such as
51 seeds of shade trees black beans, immature beans, and floaters and, when roasted, have a distinctive
52 character in the cup and high cup tasting scores.

53

54 Green bean coffee quality is a complex characteristic that depends on a series of pre-harvest factors
55 that might vary either within or among farms. Some of the pre-harvest factors that vary within
56 farms include microclimate, soil physicochemical properties, shade, age and variety of the coffee
57 tree. Factors that vary among farms include growing elevation and macroclimate, agronomic
58 practices and coffee variety (Getachew *et al.* 2022; Sarmiento-Soler *et al.* 2022). Variation in soil
59 characteristics and fertilization can be situated at both levels. Harvesting and processing conditions
60 can also be considered important factors influencing coffee quality. Elevation and shade cover
61 play an important role through temperature, availability of light and water, especially during the
62 seed ripening period (DaMatta *et al.* 2018, Sarmiento-Soler *et al.* 2022). Microclimate has a strong
63 influence on flowering, bean expansion, and ripening (Borem *et al.* 2020, Hameed *et al.* 2020).
64 Elevation is the major driving factor of climatic and edaphic factors at larger spatial scales. Shade
65 tree canopy cover, on the other hand, modulates macroclimatic trends through its effect on
66 microclimate. Cool climates (higher elevations, with at least intermediate canopy cover) had a high
67 potential to produce coffee beans possessing superior total preliminary quality, higher caffeine,

68 total chlorogenic acid (CGA) contents, and trigonelline concentrations (Worku *et al.* 2018; Tolessa
69 *et al.* 2017; paper in-press). Water deficits, on the other hand, during the coffee fruit expansion
70 and filling period caused appreciable productivity loss and decreased bean quality (Kath *et al.*
71 2020; Semedo *et al.* 2018). In terms of soil characteristics, is especially soil pH associated with
72 the acidity of coffee, body and cup cleanness. More soil nitrogen increased the caffeine content,
73 resulting in a more bitter taste of the brew (Yadessa *et al.* 2020; Clemente *et al.*, 2015). Hence,
74 assessing the importance and variations of these biophysical variables at a larger-scale across
75 elevational gradients among coffee farms and the variations among coffee trees within the farms
76 across a canopy cover gradient and their associated effect on coffee bean quality is very critical.
77 Interestingly, assessing the relative importance of within farm variability and variability among
78 farms has rarely been quantified in smallholder coffee farms. The existence of such inter-and intra-
79 farm variability could have a direct implication on input allocation, agronomic management
80 decisions and the productivity of these systems (Trevisan *et al.* 2021; Monteiro *et al.*2020; Sida *et*
81 *al.* 2020). Although several studies highlighted the effects of these variables on coffee bean quality,
82 their relative importance has not yet been quantified, and documenting such differences may help
83 to improve agronomic management decisions.

84

85 Coffee quality assessment is a key step in price setting to determine its export potential in coffee-
86 producing countries. Thus, accurate quality assessment is of major importance to many coffee
87 producers, roasters, and distributors. For each coffee bean shipment, further characterization is
88 required to verify if it attains the required quality (Dos Santos *et al.* 2016). Coffee bean quality is
89 mainly described by its physical attributes (mainly determined by bean length, diameter, and
90 hundred bean mass), raw and cup quality attributes (Cheng *et al.* 2016; Dos Santos *et al.* 2016).
91 The raw quality assessment evaluates defects that are manually separated and counted according
92 as primary and secondary defects and odor. Primary defects include full black, full sour, fungus
93 presence, foreign matter, insect damage, dried cherry, any mimic seed and soil presence whereas
94 secondary defects include partial black, partial sour, floater, immature, etc. Cup or sensory quality,
95 determine the desirability of a coffee for consumption (Tolessa *et al.* 2018; Teklu *et al.* 2011).
96 These quality attributes, such as cup cleanness, acidity, body, and flavor can be distinguished by
97 sensory organs and are assessed by professional cup tasters based on established procedures. In

98 Ethiopia, cupping score analysis enables experts to evaluate the preliminary cup assessment, which
99 is used to group coffee into different specialty categories (Ribeiro *et al.* 2021; Levate Macedo *et*
100 *al.* 2020; Okubo *et al.* 2019; Tolessa *et al.* 2016). Coffee beans graded from grade one to three are
101 grouped into specialty coffee and these three categories are further classified into different
102 specialty grades, which are called: Q1, Q2, and commercial type. In this particular study, our focus
103 is mainly on total preliminary and specialty quality and bean physical attributes.

104

105 Though standardized, methods of assessing coffee quality are prone to subjective judgments in
106 addition to being costly and time-consuming. An alternative technique is use of Near-Infrared
107 Spectroscopy (NIRS). NIRS technology-based sorting and grading systems for various aspects of
108 food quality and safety are widely used. NIRS analysis is rapid, requires limited sample
109 preparation, reduces costs of chemicals, and also multiple components can be determined on the
110 same sample from a single measurement. NIRS analyses integrated to chemometrics have been
111 proposed as an analytical methodology to characterize food (Al-Harrasi *et al.* 2020; Genisheva *et*
112 *al.* 2018), medicine (Kucharska-Ambrożej *et al.* 2020; Calvo *et al.* 2018) and coffee samples
113 (Souza *et al.* 2022; Ribeiro *et al.* 2021; Zhu *et al.* 2021). Their use has been considered to
114 distinguish coffee origin (Adnan *et al.* 2020; Giraudo *et al.* 2019; Dos Santos *et al.* 2014), assess
115 the quality of coffee beans (Tolessa *et al.* 2016; Esteban-Diez *et al.* 2004), caffeine content (Ayu
116 *et al.* 2020; Budiastira *et al.* 2018; Zhang *et al.* 2017; Pizarro *et al.* 2007) and lipids (Caporaso *et*
117 *al.* 2018) as well as to measure sugars content, roasting degrees and moisture (Levate Macedo *et*
118 *al.* 2021). NIRS, specifically for coffee, has been successfully applied to coffee analysis, including
119 determination of geographical origin (Giraudo *et al.* 2019), estimation of its chemical properties
120 (Caporaso *et al.* 2018; Ayu *et al.* 2020), roasting process monitoring (Yergenson *et al.* 2020;
121 Catelani *et al.* 2018; Bertone *et al.* 2016), adulteration detection (Chakravartula *et al.* 2022; De
122 Carvalho Couto *et al.* 2021; Correia *et al.* 2018), and sensory analysis (Ribeiro *et al.* 2011).

123

124 Nowadays, combination of multivariate calibration methods with spectroscopic data has allowed
125 the analysis of complex spectra of multi-component system. A vast range of linear and non-linear
126 computational methods is used for modeling these systems. Among them, random forest (RF) and
127 partial least square regression (PLSR) are useful when non-linear or high-dimensional

128 relationships exist in the dataset. RFs are known as a flexible approach to capture non-linear
129 relationships in high-dimensional data by learning a multitude of decision trees. Furthermore, they
130 can be used to rank the predictors according to their importance to obtain accurate predictions.
131 PLSR can cope with multidimensional data, and can eliminate multicollinearity problems by
132 generating latent variables (components) from the covariance matrix of dependent and independent
133 variables (Tyrallis *et al.* 2019).

134
135 Here we compile a database of 12 coffee trees selected across a gradient of open to dense canopy
136 cover within each of 12 farms selected across an elevational gradient from 1470 m to 2325 m
137 above sea level. This resulted in a total of 12 trees times 12 farms, that is, 144 coffee trees,
138 specifically designed to quantify among and within farm variability, and assess the relative
139 importance of the biophysical drivers of this variability. All farms were of the semi-plantation
140 coffee production type. Our response variables were the cupping scores and physical bean
141 attributes and NIRS spectra of the green coffee beans. We specifically addressed the following
142 three important questions:

- 143
- 144 i) How much green coffee bean quality variation is there among and within coffee farms?
 - 145
 - 146 ii) What is the relative importance of biophysical drivers (type and degree of canopy
147 cover, soil temperature and moisture) for green bean coffee quality?
 - 148
 - 149 iii) Can NIR spectra of green coffee beans be used to predict cupping scores?
 - 150

151
152 **2. Methodology**

153 In this section, we give an overview of the methods used in this study. In section 2.1 the data
154 acquisition is described in detail. Section 2.2 treats the preprocessing techniques that are required
155 to analyze the near infra-red spectroscopic (NIRS) data. The statistical machine learning methods
156 used to analyze the data are discussed in section 2.3.

157
158 **2.1 Data acquisition**

159 We describe all aspects related to the data acquisition. Two main types of data were acquired:
160 biophysical variables related to the selected coffee trees and NIRS data from their coffee beans.

161

162 **Study area:** The study was conducted in the Goma and Gera districts of Jimma Zone in
163 southwestern Ethiopia (7°37'48"–7°56'37"N latitude and 36°13'41"–36°39'17"E longitudes) on
164 12 coffee farms (Table 1 and Fig. 1). The region is characterized by a humid and warm subtropical
165 climate with a yearly rainfall between 1500 and 2000 mm. The main rainy season is from May to
166 September (monomodal rainfall) accounting for about 85% of the annual rainfall and coffee
167 cultivation in the region is entirely rain-fed. Differences in temperature vary throughout the year
168 with a mean monthly temperature between 13°C and 26°C. The bulk of coffee growing soils in the
169 region are classified as Eutric Nitisols, which are deep, red, and well-drained soils with a clay
170 content of more than 30% and a pH (measured in H₂O) between 4.2 and 6.2 (Muleta *et al.* 2008,
171 Dubale, 1996).

172

173 **Coffee farms selection and characterization:** The study covered agroforestry sites distributed
174 across the landscape, comprising an area of approximately 50 by 50 km. To encompass a natural
175 temperature gradient, 12 coffee farms were selected across elevational gradients ranging between
176 1470 m – 2325 m asl. All the selected coffee farms are categorized as a semi-plantation coffee
177 production system, with high anthropogenic disturbances resulting in a relatively species poor
178 canopy consisting of tree species such as *Albizia schimperiana*, *Albizia gummifera* and *Croton*
179 *macrostachyus*. Mulching and organic fertilizers are a commonly used soil fertility management
180 strategies. To avoid spatial autocorrelation, the selected farms were at least 3-4 km apart.

181

182 Within each farm, sampling was conducted in 30 x 30 m area (sampling coffee plants consistently
183 positioned inside the plantation to avoid edge effects) in which 12 individual coffee trees were
184 selected. Shade tree canopy cover was measured for each selected coffee tree. Four coffee trees
185 were sampled under each of the following canopy cover categories: light (<35% canopy cover),
186 intermediate (35-65%) and dense shade levels (>65%). Accordingly, a total of 144 coffee trees
187 were sampled from 12 coffee farms. All the measurements and data provided in this manuscript
188 (shade tree canopy cover, soil moisture content, and soil temperature) are at the individual coffee
189 tree-level with n = 144 whereas soil chemistry is at the coffee farm level.

190

191 **Biophysical variables:** The following biophysical variables were measured to describe elevation,
192 shade tree canopy cover, soil temperatures and moisture, and soil properties per individual coffee
193 trees. The variables are used to determine coffee quality (both cupping scores and NIR spectra).

194

195 **i) Elevation**

196 The elevation of each coffee farm was measured with a GPS (Garmin-60, Kansas, USA).

197

198

199 **ii) Shade tree canopy cover**

200

201 Shade tree canopy cover over each coffee tree was quantified using a convex spherical crown
202 densiometer (Forest densiometers, Model A, Bartlesville, Oklahoma, USA). The densiometer is
203 made of a small wooden box with a convex mirror consisting of a grid of squares; shade tree
204 canopy cover is then calculated as the proportion of 96 points that was intersected by vegetation
205 times 1.04. The densiometer was held at breast height and the observer's head was reflected from
206 the edge of the mirror just outside the box. The curved mirror reflects the canopy above. Above
207 the canopy of each sampled tree using a ladder all the time, two counts were recorded and their
208 mean was used.

209

210 **iii) Soil temperature, moisture and chemical characteristics**

211 a) **Soil temperature:** To quantify the temperatures in each coffee farm, soil temperatures were
212 recorded at one-hour intervals for a 10 months period (February 2020 to November 2020, i.e.,
213 period from coffee flowering to harvest) using miniature temperature sensors (type HOBO 8K
214 Pendant Temperature/Alarm Data Logger – UA-001–08, Onset Computer Corporation,
215 Bourne, MA, USA) buried in the soil at 10 cm depth and 40 cm distance from the coffee tree
216 trunk at all 144 coffee plants. We could not measure air temperatures due to theft of visible
217 devices. To ensure the best representation of temperature experienced by the coffee plant, the
218 daily minimum, mean and maximum soil temperature values were computed. From the daily
219 data, monthly mean, minimum and maximum temperatures of the period February - November
220 2020 were used for further analyses.

221

222 b) **Soil moisture (gravimetric method):** Surface mineral topsoil (0-10 cm) was sampled in the
223 dry season at the start of the coffee flowering season in February 2020 to reflect the weather-
224 independent water status of the site (rough farm ranking, independent from rainfall) using a
225 core sampler after removing the surface litter and plant debris at three locations per coffee tree
226 (10 cm away from the stem in three directions). The samples were taken during the
227 measurement of the canopy cover. These three samples were pooled into one sample for soil
228 moisture content and nutrient analysis. The mass of the fresh soil samples was recorded using
229 a balance immediately after sampling. The samples were oven-dried at 65°C for 48 hrs
230 (Robertson *et al.* 1999), after which the dry mass was recorded immediately to determine
231 gravimetric soil moisture content. Finally, the percent soil moisture was computed as (fresh
232 soil mass – dry soil mass)/dry soil mass) x 100).

233
234 c) **Soil chemical characteristics:** An oven-dried sub-sample was used for the measurements of
235 pH, soil organic carbon, total N, Olsen-P, exchangeable Ca, Mg and K (Table 1). All the soil
236 samples were dried to a constant weight at 65°C for 48 h, ground and sieved over a 2 mm
237 mesh. The pH (in H₂O) of the soil was measured using a calibrated glass electrode (model Ross
238 sure-flow 8172 BNWP, Thermo Scientific Orion, USA). Soil organic C and total N, were
239 measured using a CNS elemental analyser with a thermal conductivity detector in a (vario
240 Macro Cube, Elementar, Uberlingen, Germany). Soil total Ca, K and Mg were measured by
241 atomic absorption spectroscopy (Varian SpectrAA-220, USA) after complete destruction of
242 the soil samples with HClO₄ (65%), HNO₃ (70%) and H₂SO₄ (98%) in Teflon bombs for 4 h
243 at 150°C. Exchangeable K⁺, Ca²⁺, Mg²⁺, Na⁺ and Al³⁺ concentrations were measured by atomic
244 absorption spectroscopy after extraction in 0.1 M BaCl₂ (NEN 5738:1996).

245
246 **Coffee berry sampling and measurements:** All fully ripe, red colored coffee berries were hand-
247 picked once at peak harvest between early October and early November 2020 from each individual
248 coffee tree using local coffee bags. Berries were harvested first from lower elevation sites followed
249 by the higher elevation sites. The berries were dry processed, i.e. sun-dried (on raised beds with a
250 mesh wire) immediately after harvest (harvesting was in the morning and drying started in the
251 afternoon). The berries were returned back to the bags before sunset and stored in clean rooms (to

252 prevent any spoilage), and returned back to the raised beds in the morning until the green beans
253 attained 11.5% moisture content measured using a coffee moisture tester (mini GAC, Dickey -
254 John, USA). The berries were regularly turned to maintain uniform drying. The dried coffee berries
255 were dehusked using a coffee hulling machine (coffee huller, McKinnon, Scotland) at Jimma
256 University, cleaned and stored at room temperature in separate labeled bags for analysis.

257
258 For cup quality (60% of the total preliminary quality), green coffee bean samples were evaluated
259 for cup quality attributes by a panel of three internationally trained, experienced and certified Q-
260 grade cuppers at the Ethiopian Commodity Exchange (ECX) center based in Jimma town. Acidity,
261 body, cup cleanness and flavor were assessed in accordance with the standard method (ECX,
262 2011). This Q-grade standard method involves Q-certified cuppers, i.e., cuppers licensed by
263 Specialty Coffee Association (SCA) Coffee Quality Institute (CQI). The cuppers were trained in
264 descriptive sensory analyses in using a sensory lexicon of cup quality (Di Donfrancesco *et al.*
265 2014). Accordingly, aroma, flavor, acidity, body, uniformity, cup cleanness, overall preference,
266 aftertaste, balance and sweetness were each rated on a scale from 0 to 10. This total preliminary
267 assessment was used to classify the coffee samples into different quality grades. According to ECX
268 (2011), dry-processed coffee samples were categorized into different quality grades based on total
269 preliminary assessment and classified as: 91-100 (grade 1), 81-90 (grade 2) and 71-80 (grade 3)
270 whereas the specialty coffee achieving scores between 85-100 are classified as specialty 1 (Q₁)
271 and 80-84 is specialty 2 (Q₂), and Q₃ (commercial type) ([https://sca.coffee/research/protocols-best-](https://sca.coffee/research/protocols-best-practices)
272 [practices](https://sca.coffee/research/protocols-best-practices)).

273
274 Roasting, grinding and brew preparation was standardized. A roaster equipped with a cooling
275 system, in which air was forced through a perforated plate, capable of roasting up to 500 g of
276 coffee beans, was used for roasting the coffee beans. An amount of 100 g green beans was used
277 for each sample and the beans were put into the roasting machine with six cylinders (Probat, 4
278 Barrel Roaster, Germany) and were carefully roasted for 7-8 minutes to medium roast at
279 temperatures of 200°C. The roasted bean samples were ground to a medium level using a
280 Guatemala SB electrical grinder that was cleaned well after each sample. The medium roasted
281 coffee was tipped out into a cooling tray and allowed to cool down for 4 minutes rapidly by

282 blowing cold air through it. Then, eight gram of coffee powder was put into a 250 mL cup and 5
283 cups per sample were used. Next, 125 ml boiled water (93°C) was poured onto the ground coffee,
284 followed by stirring the content to ensure the homogeneity of the mixture. Then, the cups were
285 filled with an additional 125 mL and left to settle. After three minutes, the floaters were skimmed,
286 and the brew was ready for cup tasting. Finally, the five prepared cups were tasted by three
287 professional Q-grade cuppers operating in ECX. Each panelist gave their independent judgment
288 using a cupping form and the average score of the three cuppers was used for analysis. Total
289 preliminary quality is the sum total of raw bean quality (primary defects, secondary defects, and
290 odor) and cup quality attributes (acidity, body, flavor, and cup cleanness), whereas specialty
291 quality is the sum total of ten cup quality attributes (aroma, flavor, aftertaste, acidity, body,
292 balance, overall, cup cleanness, sweetness, and uniformity).

293

294 **Scanning of coffee bean samples with NIR spectroscopy and spectral data acquisition**

295 Approximately 50 g of each dried and grounded green coffee bean sample was placed into a glass
296 Petri dish (diameter = 2 cm, depth = 1 cm). Samples in the Petri dish were pressed gently and
297 levelled by a spatula, which was necessary as the bean powder surface ensures maximum diffuse
298 reflection and high signal-to-noise ratio. A Fourier Transform-NIR spectrometer (Tango, Brucker,
299 Belgium) was used to obtain coffee bean spectra. Green coffee samples were ground to a size
300 smaller than 5 mm and 15-20 g of each sample was used for analysis. The samples were irradiated
301 with tungsten (5V/7W) as source of near infrared light and the spectra measurements were
302 performed at room temperature. The coffee bean samples were scanned in diffuse reflectance mode
303 using a Compact NIR spectrophotometer (Tec5 Technology for spectroscopy, Germany) and the
304 reflectance was detected by an Indium Gallium Arsenide (InGaAs) diode. This generated NIR
305 spectra data consist of two lists of numbers (wavenumber and its associated reflectance). Each
306 spectrum had 1898 data points in the wavenumber range of 3952 to 11540 cm⁻¹ (867 to 2530 nm)
307 with data spacing of 4 cm⁻¹ for a total of 144 bulk coffee bean samples (Appendix Fig. S1).

308

309 **2.2 NIRS data preprocessing and analysis**

310 Preprocessing of NIR spectra is an essential component of multivariate data calibration. Its primary
311 goal is to remove unwanted information such as spectra noise, and scattering effect that are not

312 related to the variables (properties) of interest. Furthermore, we also apply outlier detection,
313 spectra trimming and optimal wavelength selection through a PCA analysis.

314
315 **Spectra preprocessing:** A variety of mathematical spectra pre-processing were tested in order to
316 improve model robustness and prediction accuracy. Several pre-processing methods such as
317 spectra normalization, de-trending (DT), standard normal variate (SNV), vector normalization,
318 spectra derivatives, multiplicative scatter correction (MSC), orthogonal signal correction (OSC)
319 and combinations of them have been investigated in several studies. Generally, with a well-tested
320 pre-processing steps, the performance of the model can be greatly improved. These mathematical
321 pre-processed methods strongly depend on a given dataset, and no universal solution could be
322 found. However, certain preprocessing techniques were selected following the best results of (Jiao
323 *et al.* 2020; Dotto *et al.* 2018; Nawar *et al.* 2017; Knox *et al.* 2015; Peng *et al.* 2014; Cambule *et*
324 *al.* 2012; Knox *et al.* 2012).

325
326 Among the extensively reviewed preprocessing techniques, Savitzky-Golay smoothing,
327 multiplicative scatter correction, and standard normal variate were found to be a better
328 preprocessing algorithms for preprocessing of our raw spectra, and all three of them were
329 examined in two models. Besides, these mathematical pre-processed algorithms were widely used
330 in reflectance spectroscopy methods in many literatures (Bian *et al.* 2021; Ren *et al.* 2021; Jiao *et*
331 *al.* 2020; Nawar *et al.* 2017). The details of these preprocessing algorithms are put in the appendix
332 word file 1.

333
334 **Outlier detection in pre-processed NIR spectra:** This was performed using the 25th and 75th
335 percentile by checking outliers in NIR spectra using the `quantile()` function from the package
336 “`ggstatsplot`”. The suspected outliers were detected using the interquartile range (IQR). Finally,
337 the `subset()` function from the same package was used to eliminate outliers. Accordingly, five rows
338 from the dataset were detected as outliers and subsequently omitted (Bello *et al.* 2020).

339
340 **Spectra trimming:** this is a procedure where wavelength ranges with high signal-to-noise ratio
341 are removed. This is to specify the wavelength regions of interest without any standard procedures

342 (Wadoux *et al.* 2021; Ng *et al.* 2018). This is because, the spectra measurements below 720 nm
343 and above 2500 nm do not contain much useful information since they are at the boundary of the
344 range recorded by the sensor. Hence, NIR spectra within a range of 720–2500 nm was retained for
345 further spectra processing, which has brought the number of data points per spectrum to $p = 1884$.

346

347 **Optimal wavelength selection:** As the complete spectra contains redundant information as
348 indicated in Appendix Fig S1, this would result in complex, unstable, and inaccurate models.
349 Hence, optimal wavelength selection is generally used to identify those wavelengths that capture
350 a large part of the information present in the spectra (Mishra *et al.* 2021; Rodriguez-Pulido *et al.*
351 2013). We performed a PCA analysis showing that the first component explained 97.8% of the
352 variability present in the whole spectra. The wavelengths corresponding to the peaks from the
353 loading plot of the first component were selected as the optimal wavelengths (Appendix Fig S3)
354 (Mishra *et al.* 2021; Rodriguez-Pulido *et al.* 2013). In this way, 87 wavelengths were selected and
355 used to infer the final model. An excel containing all the selected wavelengths using the PCA
356 loading method can be found in the Appendix table 1.

357

358 **2.3. Statistical data analysis**

359 We introduce the statistical machine learning methods that are used for data analysis.

360

361 **Variable importance:** For data visualization and to detect multicollinearity among the biophysical
362 variables, PCA was employed. The PCA provides a set of explanatory orthogonal vectors by
363 projecting similar variables in the two-dimensional space and subsequently variables closer to each
364 other indicates the high correlation. Accordingly, Tmin (minimum temp) was omitted from the
365 analysis and the remaining predictor variables were considered for variable selection procedures
366 (Janitza *et al.* 2018; Wright *et al.* 2017). Likewise, as elevation is the natural driving factor for
367 other biophysical variables, it was also omitted from the analysis. Based on this, permutation-
368 based variable importance was used to estimate the influence of a given variable in a model
369 prediction and ultimately estimate its relative importance for the coffee quality and hundred bean
370 mass. This technique assigns a score to input variables based on how useful they are at predicting
371 a target variable (Probst, 2018; Wright *et al.* 2017). Here only random forest was used for

372 classification and ranking candidate biophysical variables based on the variable importance using
373 varImp package (Probst, 2019). A higher score means that the specific variable will have a larger
374 effect on the model that is being used to predict a certain variable. By looking at the variable
375 importance, we can easily decide which variables to possibly drop because they do not contribute
376 much to the prediction process. The method for calculating permutation accuracy importance was
377 applied in R using the ranger package (Janitza *et al.* 2018; Wright *et al.* 2017).

378

379 **Permutational Multivariate Analysis of Variance (PERMANOVA):** The proportion of
380 variance explained by the individual coffee trees and coffee farms was performed using geometric
381 partitioning of variance in a multivariate data analysis technique to examine whether the variability
382 observed in physical coffee bean quality and cupping scores is between the individual coffee trees
383 or coffee farms and to quantify the proportion of variance explained by each of them. To this end,
384 a variance partitioning approach was adopted in the multivariate domain (Behrens *et al.* 2018;
385 Anderson, 2005). This relied on a PERMANOVA model, featuring the Manhattan distance matrix
386 among observations (i.e. cupping scores) as the dependent matrix and coffee trees and farms as a
387 fixed and random variable, respectively. PERMANOVA was chosen because it generates a
388 geometric partitioning and extends the analysis much broader, allowing rigorous meaningful
389 analysis of high-dimensional systems having variables with extremely non-normal or over
390 dispersed behavior. It is not restricted by distributional assumptions and can accommodate
391 heterogeneity within-group dispersions than the classical ANOVA. The model was based on 1×10^5
392 permutations and the breakdown of the variance among coffee trees and farms was carried out by
393 evaluating the marginal effect of each of them in the full model. In doing so, the share of the coffee
394 farms and coffee trees was examined through linear mixed effect models for the total preliminary
395 and specialty quality, and hundred bean mass. The estimation of the variance function, its partition
396 among the two, and how each of them affects cupping scores and hundred bean mass was then
397 performed according to the approach of Hoffman (2021). The permutational multivariate analysis
398 of variance was carried out using “variancePartition” and “vegan” package.

399

400 **Random forest and partial least square regression to predict coffee quality based on NIR**
401 **spectra data:** A two-dimensional data matrix consisting of pre-processed spectra (as independent

402 variables) and measured cupping scores (total preliminary and specialty quality) as dependent
403 variable was created from the 139 coffee bean samples (as 5 of the rows were outliers and
404 subsequently removed). For the purpose of coffee quality prediction random forests (RF) and
405 partial least square regression (PLSR) were tested.

406
407 Both PLSR and RF models were chosen because they are more powerful than the conventional
408 regression models for modeling complex and non-linear data in a high-dimensional and
409 hierarchical fashion. PLSR can cope with multidimensional data, and can eliminate
410 multicollinearity problems by generating latent variables (components) from the covariance matrix
411 of dependent and independent variables. Hence, PLSR is recommended as one of the best
412 performing calibration techniques for spectral data (Kuang *et al.* 2015). RF, on the other hand,
413 uses an ensemble of a large number of decision trees by offering sufficient accuracy, simple
414 implementation, and high robustness (Tyralis *et al.* 2019). The algorithm is a model ensemble
415 method constructed based on combining several decisions by the regression and classification
416 trees. Two key parameters need to be taken into account: one is the number of the decision trees
417 and the other is the number of sampled variables for building a decision tree. RF has the capability
418 of ranking the importance of variables by their importance (Janitza *et al.* 2018). The method could
419 be briefly summarized in three steps: (1) the Bagging method to generate T subsets of training data
420 randomly; (2) each training sample is employed to generate the corresponding decision trees
421 randomly choosing m attributes from M attributes as the split attributes set of the current node
422 prior to select attributes on each non-leaf node, and split the node in the best split way among the
423 M attributes; (3) each tree grows sufficiently without pruning, and was used to test the
424 corresponding category from the test set. Finally, the majority vote of the decision trees was used
425 to make an ensemble classification decision (He *et al.* 2022; Khan *et al.* 2022; Asadi *et al.* 2021;
426 Ao *et al.* 2019).

427
428 **Model validation and evaluation:** To validate the RF and PLSR for coffee quality prediction, a
429 leave-one-out cross-validation (LOOCV) procedure was performed. In each run, one sample was
430 left out to test the models while the other samples are used to train and calibrate the models.
431 Training and calibration involved a randomly split into calibration (n=112) and prediction (n=27)

432 samples. For the PLSR the number of latent variables were optimized during this calibration step.
433 For the RF, we optimized the number of sampled variables (over a range of 2 to 87) and the
434 minimum size of terminal nodes (over the values 1,5 or 10) while the number of trees was held
435 fixed at 500 (Tridawati *et al.* 2020; Wadoux *et al.* 2019; Freeman *et al.* 2016). The process of
436 LOOCV was repeated until every sample is left out once. The implementation was performed with
437 the R environment for statistical computing using the packages ‘caret’ and ‘randomForest’ (Kuhn
438 and Johnson, 2013).

439
440 In both RF and PLSR models, four statistical metrics: correlation coefficient (r), coefficient of
441 determination (R^2), root mean square error (RMSE), and residual predictive deviation (RPD) were
442 used to evaluate the predictive performance of the models according to the classification criteria
443 of Viscarra *et al.* (2009). The coefficient of determination (R^2) reflects the percentage of variance
444 in the response variable that is accounted for by the explanatory variables. An R^2 value between
445 0.5-0.65 indicates that more than half of the variance in the response variable is accounted for by
446 the explanatory variable. R^2 value in the range of 0.66-0.81 indicates approximate quantitative
447 predictions whereas the R^2 value in the range of 0.82-0.9 reveals a good prediction. Calibration
448 models possessing an R^2 value above 0.91 are considered excellent (Nakagawa *et al.* 2017). RMSE
449 allows to measure how far the predicted values deviate from the observed values in a regression
450 analysis. The larger the difference, the larger the gap between the predicted and observed values.
451 The smallest RMSE value is usually related to the optimal calibration model and the better a model
452 is able to fit the data. RPD takes both the prediction error and the variation of observed values into
453 account, hence providing a metric of model validity that is more objective than RMSE and more
454 easily comparable across model validation studies. The greater the RPD, the better the model's
455 predictive capability (Nakagawa *et al.* 2017; Kapper *et al.* 2012).

456
457 RPD is defined as the standard deviation of the measured value divided by the RMSE of the
458 predicted values (Kapper *et al.* 2012; Guy *et al.* 2011). It is calculated as follows:

459

$$RPD = \frac{SD}{RMSE}$$

460 where SD is the standard deviation of the measured value and RMSE is the standard error of
461 prediction. In general, when $RPD \geq 2$, it indicates that the model works well and can be used for

462 quantitative analysis and evaluation (Kapper *et al.* 2012). Models estimations were computed
463 using PLSR and RF along with four statistical metrics: correlation coefficient (r), coefficient of
464 determination (R^2), root mean square error (RMSE), and residual predictive deviation (RPD)
465 according to the classification criteria of Viscarra *et al.* (2009). Generally, a good model prediction
466 corresponds to high R^2 , r , and RPD, and low RMSE values. Finally, scatter plots showing the
467 relationship between the spectra data and cupping scores were generated using the best model. For
468 all statistical procedures, R-4.1.2 software (R Core Team, 2022) was used.

469

470 **3. Results**

471 Descriptive statistical results of soil chemistry, moisture and temperature in the twelve coffee
472 farms are shown in Table 1. The coefficients of variation (CV) of soil chemistry, moisture and
473 temperature showed that among all the measured soil chemical variables, Olsen-P had the highest
474 CV (63.7%), particularly in the coffee farm with elevation 2325 m asl followed by soil
475 exchangeable K (62.7%) at coffee farms situated at 2027 m asl elevation. Likewise, Olsen-P had
476 the next highest CV (60.0%) at the coffee farm situated at the elevation of 1774 m asl, as compared
477 to other measured soil parameters. In contrast, soil C at the coffee farm situated at the elevation of
478 1650 m asl had the lowest CV (0.1%), showing that this soil chemical variable is more homogenous
479 than other soil chemical variables in the study sites. CV values of soil moisture content showed
480 inconsistent values across elevational gradients.

481

482 **3.1. The relationship between elevation and canopy cover on coffee quality attributes**

483 Elevation significantly ($p < 0.05$) affected all the three coffee bean quality attributes (total
484 preliminary quality, specialty quality, and hundred bean mass). An interaction effect of elevation
485 and canopy cover significantly influenced total preliminary quality and hundred bean mass but not
486 the specialty quality (Fig. 2 and Appendix table S1). A clear relationship of hundred bean mass
487 with elevation and shade canopy cover was observed, confirming that a greater hundred bean mass
488 was produced in response to increasing elevations, at intermediate and dense shade levels (Fig. 2
489 and Appendix table S1). Under conditions of light shade levels (10-35%) and intermediate shade
490 levels (35-65%), hundred bean mass increased with elevation (Fig. 2 and Appendix table S1).
491 However, a decreasing trend in hundred bean mass was observed under dense shade (>65% shade

492 level) when elevation keep increasing. Dense shaded conditions and the low temperatures at higher
493 elevations may not improve the growth potential and quality of the beans, however, dense shaded
494 environments at warmer environments of the low elevations show an increased trend in bean mass.
495 Our study confirmed that higher elevations with cooler climates (>1900 m) with intermediate
496 shade cover (35-65%) showed a higher potential to produce green coffee beans having superior
497 total preliminary quality (Fig. 2 and Appendix table S1). The results further support that coffee
498 cup quality attributes are more sensitive to temperature changes than to other farm management
499 practices, possibly due to the fact that cup quality attributes such as flavor, taste, aroma and body
500 are temperature-dependent. Coffee beans from higher elevations had a greater specialty quality as
501 compared to coffee beans grow in warmer climates (Fig. 2 and Appendix table S1). Shade cover
502 affected the total preliminary quality and it had no effect on specialty quality. This confirms that
503 shade drives the physical bean quality (raw value) much more than the sensory attributes. In other
504 words, poor management of shade at a given elevation will have a negative impact on the potential
505 to produce qualitative coffee.

506

507 **3.2. Quantifying the proportion of variance explained by coffee trees and coffee farms**

508

509 The permutational multivariate analysis of variance (PERMANOVA) depicted significant
510 differences in coffee quality among coffee farms and between the individual coffee trees (model
511 residuals). In terms of the breakdown of the total variance in total preliminary quality, a linear
512 mixed model depicted a substantial contribution of individual coffee trees (73%) while only 17%
513 was explained by the coffee farm (Fig. 3a). In the specialty quality, the model depicted large
514 contribution of the individual coffee trees (96%), and only 4% by the coffee farms (Fig. 3a).
515 Similarly, for hundred bean mass, the model depicted a substantial contribution of individual
516 coffee tree (76.6%) and coffee farms (23.4%) (Fig. 3a).

517

518 Meanwhile, the biophysical variables contributed differently for the coffee quality attributes.
519 Canopy cover contributed 9.6%, 4.2% and 22.4% for total preliminary quality, specialty quality
520 and hundred bean mass, respectively. Soil moisture contributed 0.1%, 0.2% and 2.7% for total
521 preliminary quality, specialty quality and hundred bean mass, respectively. Tmax (max
522 temperature) contributed 2.1%, 2.3% and 0.2% for total preliminary quality, specialty quality and

523 hundred bean mass, respectively. Tmean (mean temperature) contributed 0.1% each for total
524 preliminary quality, specialty quality and hundred bean mass, respectively (Fig. 3b).

525 **3.3. Establishing a relationship between NIRS and cup quality**

526

527 The outcomes of the examined models in the estimation of coffee quality and their performance
528 assessment using different statistical metrics are presented in Table 2. Both RF and PLSR models
529 were tested, optimized and compared to each other. A good model prediction corresponds to high
530 R^2 , r , and RPD, and low RMSE values. As can be seen from table 2, the performance of the models
531 without the spectra preprocessing was low, and spectral data preprocessing could significantly
532 improve the performance of the two models. Therefore, preprocessing of the raw spectral data is
533 an important first stage before any regression model is established as it improves the prediction.
534 In addition, compared with the single preprocessing method, the combination of different
535 preprocessing methods can greatly improve the performance of the model. Accordingly, the results
536 suggest that the RF model has a better predictive power as compared to PLSR for both training
537 and testing datasets at a specified preprocessing algorithms for both total preliminary and specialty
538 quality as indicated in (Table 2 and Fig. 4). Moreover, RF model showed a higher R^2 and lower
539 RMSE values as compared to PLSR in the estimation of total preliminary and specialty quality,
540 which demonstrated that it has satisfactory estimative capability in coffee quality assessment. RPD
541 in total preliminary and specialty quality was also found to be superior in RF model (Table 2).

542

543 **3.4. Quantifying the relative importance of biophysical variables on measured and predicted** 544 **coffee quality**

545

546 By applying random forest models, the main important biophysical variables influencing coffee
547 quality were identified. The measured relative importance of the investigated variables derived
548 from the RF model as shown in Fig. 5 differed among the different coffee cupping scores. As
549 presented in the figure, the order of importance of the variables to total preliminary and specialty
550 quality and NIRS is: canopycover>soilmoisture>Tmean>Tmax, in which the first three explained
551 31.8 - 40%, 27.4 – 33.3% and 19.2%, respectively, of the variation in the data. On the other hand,
552 the order of importance of the variables to hundred bean mass is:

553 canopycover>Tmean>soilmoisture>Tmax, in which the first three explained 32.4%, 29.3% and
554 22.2%, respectively, of the variation in the data (Fig. 5). Similarly, the order of importance of the
555 variables to NIRS is as follows: canopycover>soilmoisture>Tmean>Tmax, in which the first three
556 of them contributed 37.7%, 28.2% and 22.4%, respectively for the variations. Hence, the
557 biophysical variables affecting coffee cupping (total preliminary and specialty quality) appeared
558 to be the same for the NIR spectra.

559

560 **4. Discussion**

561

562 **4.1. Intra-farm variability is larger than the inter-farm variability**

563

564 Our findings indicate that elevation is the key biophysical variable influencing all the measured
565 coffee bean quality attributes (hundred bean mass, total preliminary and specialty quality) at the
566 farm level (Fig 2 and Appendix table 1) while canopy cover was the most important biophysical
567 variable driving the coffee bean quality attributes and NIRS at the plant level (Fig 3).

568

569 Most importantly, the results show the existence of high variability between coffee plants within
570 a farm, as evidenced from the variance partitioning procedures in permutational multivariate
571 analysis of variance in a linear mixed model. The magnitude of variability observed within a coffee
572 farm is far larger than the variability among coffee farms in terms of the measured coffee bean
573 quality attributes. The order of importance of the variables to total preliminary and specialty
574 quality was found to be in order of canopycover>soilmoisture>Tmean>Tmax. This means that
575 conditions varying at the coffee plant-level might be of greater importance for influencing hundred
576 bean mass and cupping scores when considering the farm-level as a whole. The potential
577 explanations for this huge intra-farm variability could be due to various reasons: variation in
578 genetic structure of the coffee plants (there is definitely an inherent variation in growth rate among
579 coffee plants due to the variation in resource use (for instance, nutrient capture, transport and
580 utilization efficiency, water use efficiency, light use efficiency), individual leaf trait variability
581 like SLA due to the variation in genetic structure of the coffee plants, disease sensitivity of the
582 individual coffee plants); the way how the coffee plants were obtained (if not all, most coffee
583 plants of the smallholders are reared from seeds by natural means, and this could be a potential

584 explanation for the huge intra-farm variability). Variation will become high in the case of sexually
585 reproduced plants especially if it is reared by natural selection); poor and inconsistent farm
586 management practices within a farm (for instance, no definite spacing between plants and rows in
587 smallholder coffee farms unlike that of the plantations); age of the coffee plants (although an effort
588 was made to consider coffee trees aged between 6 and 10 years old, we still believe that age
589 matters).

590

591 Although significant variability has been documented in many African smallholder production
592 systems, agronomic research for development generally ignores such variability in the decision-
593 making process and programs (Oyinbo *et al.* 2019; MacCarthy *et al.* 2018). Incorporating this
594 variability in agronomic decision-making to minimize its effect requires systematic quantification
595 of the variability. However, quantifying intra-farm variability has been a challenge so far and
596 operationalizing this variability in agronomic decision-making is even more challenging (Sida *et*
597 *al.* 2021; Trevisan *et al.* 2021; Van Loon *et al.* 2019).

598

599 Most of the observed variation (more than 75%) is due to unmeasured variation or residuals. This
600 implies that a large proportion of the variation in coffee cupping scores is to be explained by other
601 plant-level factors such as the specific nutrient levels, coffee tree pruning, fruit thinning, rate,
602 method and timing of fertilizer application, age of the coffee trees, shade tree species, cultivar
603 characteristics, and disease sensitivity of the individual coffee trees. Besides, the variance
604 partitioning procedures have shown that the total preliminary and specialty quality, and hundred
605 bean mass, were driven by the shade tree canopy cover (30-46% variability), which reflects that
606 canopy cover at the coffee plant-level is more important for explaining variation for green bean
607 quality. Hence, based up on the local conditions and the requirements, smallholder coffee farmers
608 can manage their shade tree canopy cover to optimize their coffee quality.

609

610 In addition, Olsen-P had the highest CV (63.7%), particularly in the coffee farm with elevation
611 2325 m asl followed by soil exchangeable K (62.7%) at coffee farms situated at 2027 m asl
612 elevation. It is interesting to notice that, unexpectedly, soil available P content (in this case, Olsen-
613 P) was higher at higher elevations compared to the lower and mid-elevations. The most likely

614 reason for this could be differences in local environmental conditions mainly soil characteristics
615 such as weathering and/or litter quality. We still believe that more samples would have improved
616 the accuracy of the fluctuation in P concentrations. Consequently, more data would be necessary
617 to test this hypothesis.

618

619 **4.2. Importance of the biophysical variables for coffee cupping scores and NIRS**

620

621 The results of the random forest model indicated that canopy cover, soil moisture and mean soil
622 temperature were identified as the key variables affecting total preliminary and specialty quality,
623 and hundred bean mass at a coffee-tree level. Large number of studies have shown that shade
624 percentage, soil moisture and temperature affect coffee cupping scores. At the local scale, canopy
625 cover is the main determinant of microclimate temperature and radiation. Shade canopy cover
626 provides a means to keep coffee plants closer to their ideal temperature ranges (18°C-21°C) and
627 prevent damage from extreme minimum and maximum temperatures and drought (Nesper *et al.*
628 2017, Somporn *et al.* 2012). Numerous studies have shown that there was a significant positive
629 correlation between coffee quality and shade tree canopy cover as well as soil moisture and
630 significant negative correlation between coffee quality and temperature (Bosselmann *et al.* 2009;
631 Avelino *et al.* 2007; Leonel and Philippe, 2007). A decline in soil temperature were recorded with
632 elevation, implying that the spatial distribution of soil temperatures is controlled mainly by
633 elevation (Navarro-Serrano *et al.* 2020). Although soil temperatures can be affected by the
634 interaction of multiple local factors such as shade canopy cover, mulching and irrigation, elevation
635 was found to be the main driving variable for the changes in soil temperatures (Getachew *et al.*
636 2022; paper in-press; Onwuka and Mang, 2018; Barman *et al.* 2017).

637

638 Likewise, DaMatta *et al.* (2018) reported that coffee plants tolerated higher temperatures when
639 ample water was supplied. A study from Southwest Ethiopia demonstrated that coffee trees grown
640 under open shade conditions produced beans of lower acidity, body and flavor as compared to the
641 coffee plants grown under dense shade (Bote, 2016). On the other hand, higher bean size and mass
642 were obtained when shade canopy cover increased. Shade promotes slower and more balanced
643 fruit maturation by the mother plant, thus yielding a better-quality product than unshaded coffee
644 plants (Barbosa *et al.* 2012; Geromel *et al.* 2008; Leonel and Philippe, 2007). These previous

645 findings further support that coffee cup qualities are more sensitive to temperature changes,
646 possibly due to the fact that formation of biochemical precursor molecules responsible for cup
647 quality attributes such as flavor, taste, aroma and body are temperature dependent. Meanwhile,
648 Bertrand *et al.* 2012 demonstrated that mean soil temperature during coffee bean development
649 influenced acidity, fruity character and flavor. Silva *et al.* (2005) reported that temperature was
650 likely the most important factor to bring variations in coffee cup quality from the southwest region
651 of Ethiopia. Our results thus corroborate that the physical attributes and cupping scores are more
652 temperature driven. Altitude and shade cover management is therefore important to enhance the
653 potential to bring good coffee beans to the market.

654 **4.3. Comparison of the two models for the quantitative prediction of coffee cupping scores**

655 Based on two of the tested models (RF and PLSR), the effects of different preprocessing methods
656 were examined. The pre-processing of spectral data can remove the influence of irrelevant
657 information on our spectra and ultimately improved the robustness and accuracy of the models. As
658 can be seen from table 2, RF and PLSR models produced different outputs when different
659 preprocessing methods were used separately or in combination. When the spectra were completely
660 not preprocessed, R^2 was only 0.56 and 0.52, and RMSE was 0.87 and 1.92 in PLSR and RF,
661 respectively in specialty quality. After application of Savitzky-Golay smoothing, multiplicative
662 scatter correction, and standard normal variate, the R^2 was raised to 0.87, while RMSE was reduced
663 to 0.26 when RF was used. Therefore, the performance of the RF model without preprocessing
664 was obviously low, and spectral data preprocessing could significantly improve the performance
665 of the model. In addition, compared with the single preprocessing method, the combination of
666 different preprocessing methods was of great help to the performance.

669 Referring to Table 2 again, the RMSE of the RF for total preliminary quality was 0.55, which
670 represented the lowest error of prediction when Savitzky-Golay smoothing, multiplicative scatter
671 correction and standard normal variate preprocessing methods were applied. The R^2 is 0.83,
672 indicating that the RF model can predict the data reasonably better (Barea-Sepulveda *et al.* 2022,
673 Anderson *et al.* 2020; Zhang *et al.* 2020; Ghasemi and Tavakoli, 2013). Most importantly, the
674 RPD of random forest was 3.87, whereas the RPD of PLSR was 1.41 for the total preliminary
675 quality when the same preprocessing methods were utilized. The RF thus performed well (Barea-
676

677 Sepulveda *et al.* 2022, Anderson *et al.* 2020). Given these findings, a simultaneous application of
678 spectral preprocessing methods (Savitzky-Golay smoothing, multiplicative scatter correction and
679 standard normal variate) in conjunction with the RF model better predicted the coffee cup quality.

680
681 Several studies have reported that the performance of different tree-based models including RF
682 can vary from study to study, thus there is no general best modelling technique for predicting
683 coffee quality (NS Akbar *et al.* 2020; Martinez-Santos *et al.* 2021). Moreover, it has been
684 suggested that the predictive power of the modelled output is also the result of the research design,
685 preprocessing methods and input variables (Vargas and Hanandeh, 2021; Naccarato *et al.* 2016;
686 Aertsen *et al.* 2010). Overall, near-infrared spectroscopic based predictions of green bean quality
687 can be utilized to complement cupping evaluations conducted by humans, and most importantly,
688 to increase the throughput of the cupping evaluations.

689

690 **Limitations and way forward**

691 Our results show the existence of high variability among coffee plants within farms, which can be
692 as high as 73%. Although we have quantified the magnitude and distribution of the inter-and intra-
693 farm variability in smallholder coffee farms, a couple of questions remain unaddressed in our
694 study. We were limited to disentangle the drivers of some relationships in this work because of
695 data limitations such as coffee cultivar characteristics (as there is definitely an inherent variation
696 in growth rate among coffee plants due to the variation in resource use), disease sensitivity of the
697 individual coffee trees, limited soil moisture data, plant nutrient levels, etc. Most importantly, our
698 study is also a relatively a short-term study and this again calls for caution for generalizing our
699 results and long-term investigations are necessary.

700 **Conclusion**

701 The main contribution of this work is the assessment of the spatial variability of coffee quality in
702 response to different biophysical drivers. Our study showed the existence of large within farm
703 variability, indicating that conditions varying at the coffee plant-level are of importance for
704 improving the physical attributes and cupping scores of green coffee beans, and hence
705 documenting such differences may help to improve agronomic decision-making processes.
706 However, quantifying the factors responsible for the large within farm variability is much more

707 challenging than identifying and measuring among-farm variability. Understanding, quantifying,
708 and managing within farm variability is crucial to improve to improve nutrient use efficiency,
709 water availability, pruning, pest and disease control, etc. Meanwhile, the overall biophysical
710 variables responsible for the coffee cupping scores and NIRS were identified and quantified, which
711 are fundamental to improving coffee quality. Overall, elevation was the key variable driving
712 biophysical variable influencing all the measured coffee bean quality attributes (hundred bean
713 mass, total preliminary and specialty quality) at the farm level while canopy cover was the most
714 important biophysical variable driving the coffee bean quality attributes and NIRS at the shrub
715 level. Accordingly, canopy cover appeared to be the main controlling variable for the variation in
716 total preliminary and specialty quality, and hundred bean mass, followed by soil moisture and soil
717 temperatures. On the other hand, NIRS was confirmed to be a good approach in estimating the
718 cupping scores. However, the developed NIRS models need to be tested further on data from other
719 Ethiopian regions to ensure the models' stability and accuracy.

720 **Acknowledgements**

721
722 This study has been supported by the Belgian Development Cooperation (NASCERE program)
723 and Ethiopian Ministry of Science and Higher Education (MoSHE). The authors are profoundly
724 grateful for the Ethiopian government for this support. We would like to thank coffee owners for
725 allowing us to work in their coffee plots and the local and regional administration for providing
726 the permits to work in the coffee farms. We also thank Beyene Zewdie for establishing the coffee
727 plots. Lastly, we are very grateful to the ECX for their support in evaluating coffee cup quality.

728

729 **Code availability:** Not applicable.

730 **Authors' contribution:** Conceptualization, M.G., P.B., K.T., P.D.F. and K.V.; methodology,
731 M.G., K.T., K.V., K.H., A.T., B.A., P.B. and P.D.F.; formal analysis, M.G. and S.L.; investigation,
732 M.G. and K.T., writing—original draft, M.G. and P.D.F.; writing—review and editing, P.B., K.H.,
733 A.T., and K.V.; visualization, M.G., and B.A.; Funding Acquisition, P.B. and P.D.F.; Supervision,
734 K.V., P.B. and P.D.F.

735

736 **Funding:** The study has been supported by the NASCERE program of the Ethiopian Ministry of
737 Science and Higher Education and Global Minds from Ghent University. The authors are
738 profoundly grateful for the Ethiopian government for this support. P.D.F. received funding from
739 the European Research Council (ERC) under the European Union's Horizon 2020 research and
740 innovation programme (ERC Starting Grant FORMICA).

741

742 **Data availability:** The datasets and R-code analyzed during the current study will be made
743 available after publication via an online repository such as figshare.

744 **Declarations**

745 **Conflict of interest:** The authors declare no potential conflict of interest.

746 **Ethics approval:** Not applicable.

747 **Consent to participate:** Not applicable.

748 **Consent for publication:** Not applicable.

749

750 **References**

751

752 Adnan, A., Naumann, M., Morlein, D. and Pawelzik, E., 2020. Reliable discrimination of green
753 coffee beans species: A comparison of UV-Vis-based determination of caffeine and
754 chlorogenic acid with non-targeted near-infrared spectroscopy. *Foods*. 9(6):788.

755

756 Aertsen, W., Kint, V., Van Orshoven, J., Ozkan, K. and Muys, B., 2010. Comparison and ranking
757 of different modelling techniques for prediction of site index in Mediterranean mountain
758 forests. *Ecological modelling*, 221(8):1119-1130.

759

760 Al-Harrasi, A., Rehman, N.U., Mabood, F., Albroumi, M., Ali, L., Hussain, J., Hussain, H., Csuk,
761 R., Khan, A.L., Alam, T. and Alameri, S., 2017. Application of NIRS coupled with PLS
762 regression as a rapid, non-destructive alternative method for quantification of KBA in
763 *Boswellia sacra*. *Spectrochimica Acta Part A: Molecular and Biomolecular*
764 *Spectroscopy*, 184(5):277-285.

765

766 Anderson, M.J., 2005. Permutational multivariate analysis of variance. Department of Statistics,
767 University of Auckland, Auckland, 26:32-46.

768

- 769 Anderson, N.T., Walsh, K.B., Subedi, P.P. and Hayes, C.H., 2020. Achieving robustness across
770 season, location and cultivar for a NIRS model for intact mango fruit dry matter
771 content. *Postharvest Biology and Technology*, 168:111202.
772
- 773 Ao, Y., Li, H., Zhu, L., Ali, S. and Yang, Z., 2019. The linear random forest algorithm and its
774 advantages in machine learning assisted logging regression modeling. *Journal of*
775 *Petroleum Science and Engineering*, 174:776-789.
776
- 777 Asadi, S., Roshan, S. and Kattan, M.W., 2021. Random forest swarm optimization-based for heart
778 diseases diagnosis. *Journal of Biomedical Informatics*, 115:103690.
779
- 780 Avelino J, Barboza B, Davrieux F and Guyot B. 2007. Shade effects on sensory and chemical
781 characteristics of coffee from very high-altitude plantations in Costa Rica. In *Second*
782 *International Symposium on Multi-strata Agroforestry Systems with Perennial Crops.*
783 *September 17-21. Turrialba, Costa Rica. Oral and poster presentations. Turrialba: CATIE*
784
- 785 Ayu, P.C., Budiastara, I.W. and Rindang, A., 2020, February. NIR spectroscopy application for
786 determination of the caffeine content of Arabica green bean coffee. In *IOP Conference*
787 *Series: Earth and Environmental Science*. 454(1):012049).
788
- 789 Barbosa, J.N., Borem, F.M., Cirillo, M.A., Malta, M.R., Alvarenga, A.A. and Alves, H.M.R. 2012.
790 Coffee quality and its interactions with environmental factors in Minas Gerais, Brazil. *J*
791 *Agr Sci*. 4(5):181.
792
- 793 Barea-Sepulveda, M., Ferreiro-González, M., Calle, J.L.P., Barbero, G.F., Ayuso, J. and Palma,
794 M., 2022. Comparison of different processing approaches by SVM and RF on HS-MS
795 eNose and NIR Spectrometry data for the discrimination of gasoline
796 samples. *Microchemical Journal*. 172:106893.
797
- 798 Barman, D., Kundu, D.K., Pal, S., Chakraborty, A.K., Jha, A.K., Mazumdar, S.P., Saha, R. and
799 Bhattacharyya, P., 2017. Soil temperature prediction from air temperature for alluvial soils
800 in lower Indo-Gangetic plain. *Int Agrophys*. 31(1).
801
- 802 Behrens, T., Schmidt, K., Viscarra Rossel, R.A., Gries, P., Scholten, T. and MacMillan, R.A.,
803 2018. Spatial modelling with Euclidean distance fields and machine learning. *European*
804 *journal of soil science*, 69(5):757-770.
805
- 806 Bello, M., Napoles, G., Morera, R., Vanhoof, K. and Bello, R., 2020, October. Outliers detection
807 in multi-label datasets. In *Mexican International Conference on Artificial Intelligence*.
808 12468:65-75.
809
- 810 Bertone, E., Venturello, A., Giraud, A., Pellegrino, G. and Geobaldo, F.J.F.C., 2016.
811 Simultaneous determination by NIR spectroscopy of the roasting degree and
812 Arabica/Robusta ratio in roasted and ground coffee. *Food Control*. 59:683-689.
813

- 814 Bertrand, B., Boulanger, R., Dussert, S., Ribeyre, F., Berthiot, L., Descroix, F. and Joet, T. 2012.
815 Climatic factors directly impact the volatile organic compound fingerprint in green Arabica
816 coffee bean as well as coffee beverage quality. *Food Chem.* 135(4):2575-2583.
817
- 818 Bertrand, B., Etienne, H., Lashermes, P., Guyot, B. and Davrieux, F., 2005. Can near-infrared
819 reflectance of green coffee be used to detect introgression in *Coffea arabica*
820 cultivars?. *Journal of the Science of Food and Agriculture*, 85(6):955-962.
821
- 822 Bian, X., Wang, K., Tan, E., Diwu, P., Zhang, F. and Guo, Y., 2020. A selective ensemble
823 preprocessing strategy for near-infrared spectral quantitative analysis of complex
824 samples. *Chemometrics and Intelligent Laboratory Systems*, 197:103916.
825
- 826 Borém, F.M., Cirillo, M.Â., de Carvalho Alves, A.P., dos Santos, C.M., Liska, G.R., Ramos, M.F.
827 and de Lima, R.R., 2020. Coffee sensory quality study based on spatial distribution in the
828 Mantiqueira mountain region of Brazil. *Journal of Sensory Studies*, 35(2):e12552.
829
- 830 Bosselmann, A.S., Dons, K., Oberthur, T., Olsen, C.S., Ræbild, A. and Usma, H., 2009. The
831 influence of shade trees on coffee quality in small holder coffee agroforestry systems in
832 Southern Colombia. *Agriculture, ecosystems & environment*, 129(1-3):253-260.
833
- 834 Bote, A., 2016. Examining growth, yield and bean quality of Ethiopian coffee trees: towards
835 optimizing resources and tree management (Doctoral dissertation, Wageningen
836 University).
837
- 838 Budiastira, I.W., Widyotomo, S. and Ayu, P.C., 2018, May. Prediction of caffeine content in java
839 preanger coffee beans by NIR spectroscopy using PLS and MLR method. In IOP
840 conference series: earth and environmental science. 147(1):e012004).
841
- 842 Buendia Garcia, J., Gornay, J., Lacoue-Negre, M., Mas Garcia, S., Er-Rmyly, J., Bendoula, R. and
843 Roger, J.M., 2022. A novel methodology for determining effectiveness of preprocessing
844 methods in reducing undesired spectral variability in near infrared spectra. *Journal of Near*
845 *Infrared Spectroscopy*, 30(2):74-88.
846
- 847 Byrareddy, V., Kouadio, L., Mushtaq, S., Kath, J. and Stone, R., 2021. Coping with drought:
848 Lessons learned from robusta coffee growers in Vietnam. *Climate Services*, 22:100229.
849
- 850 Calvo, N.L., Maggio, R.M. and Kaufman, T.S., 2018. Characterization of pharmaceutically
851 relevant materials at the solid state employing chemometrics methods. *Journal of*
852 *pharmaceutical and biomedical analysis*, 147:538-564.
853
- 854 Cambule, A.H., Rossiter, D.G., Stoorvogel, J.J. and Smaling, E.M.A., 2012. Building a near
855 infrared spectral library for soil organic carbon estimation in the Limpopo National Park,
856 Mozambique. *Geoderma*, 183:41-48.
857

- 858 Caporaso, N., Whitworth, M.B., Grebby, S. and Fisk, I.D., 2018. Rapid prediction of single green
859 coffee bean moisture and lipid content by hyperspectral imaging. *Journal of food*
860 *engineering*, 227:18-29.
861
- 862 Catelani, T.A., Santos, J.R., Páscoa, R.N., Pezza, L., Pezza, H.R. and Lopes, J.A., 2018. Real-time
863 monitoring of a coffee roasting process with near infrared spectroscopy using multivariate
864 statistical analysis: A feasibility study. *Talanta*, 179:292-299.
865
- 866 Chai, T. and Draxler, R.R., 2014. Root mean square error (RMSE) or mean absolute error (MAE):
867 Arguments against avoiding RMSE in the literature. *Geoscientific model*
868 *development*, 7(3):1247-1250.
869
- 870 Chakravartula, S.S.N., Moscetti, R., Bedini, G., Nardella, M. and Massantini, R., 2022. Use of
871 convolutional neural network (CNN) combined with FT-NIR spectroscopy to predict food
872 adulteration: A case study on coffee. *Food Control*:108816.
873
- 874 Cheng, B., Furtado, A., Smyth, H.E. and Henry, R.J. 2016. Influence of genotype and environment
875 on coffee quality. *Trends Food Sci Tech.* 20-30.
876
- 877 Clemente, J.M., Martinez, H.E.P., Alves, L.C., Finger, F.L. and Cecon, P.R., 2015. Effects of
878 nitrogen and potassium on the chemical composition of coffee beans and on beverage
879 quality. *Acta Sci-Agron.* 297-305.
880
- 881 Correia, R.M., Tosato, F., Domingos, E., Rodrigues, R.R., Aquino, L.F.M., Filgueiras, P.R.,
882 Lacerda Jr, V. and Romao, W., 2018. Portable near infrared spectroscopy applied to quality
883 control of Brazilian coffee. *Talanta*, 176:59-68.
884
- 885 Cozzolino, D., Kwiatkowski, M.J., Parker, M., Cynkar, W.U., Dambergs, R.G., Gishen, M. and
886 Herderich, M.J., 2004. Prediction of phenolic compounds in red wine fermentations by
887 visible and near infrared spectroscopy. *Analytica Chimica Acta*, 513(1):73-80.
888
- 889 DaMatta, F.M., Avila, R.T., Cardoso, A.A., Martins, S.C. and Ramalho, J.C., 2018. Physiological
890 and agronomic performance of the coffee crop in the context of climate change and global
891 warming: A review. *J Agr food Chem*: 5264-5274.
892
- 893 Davis, A.P., Chadburn, H., Moat, J., O'Sullivan, R., Hargreaves, S. and Nic Lughadha, E., 2019.
894 High extinction risk for wild coffee species and implications for coffee sector
895 sustainability. *Science advances*, 5(1):3473.
896
- 897 De Carvalho Couto, C., Freitas-Silva, O., Morais Oliveira, E.M., Sousa, C. and Casal, S., 2021.
898 Near-Infrared Spectroscopy Applied to the Detection of Multiple Adulterants in Roasted
899 and Ground Arabica Coffee. *Foods*, 11(1):61.
900

- 901 De Sousa, M.M., Carvalho, F.M. and Pereira, R.G., 2020. Colour and shape of design elements of
902 the packaging labels influence consumer expectations and hedonic judgments of specialty
903 coffee. *Food Quality and Preference*, 83:103902.
904
- 905 Di Donfrancesco, B., Gutierrez Guzman, N. and Chambers IV, E., 2014. Comparison of Results
906 from Cupping and Descriptive Sensory Analysis of C olombian Brewed Coffee. *Journal of*
907 *Sensory Studies*. 29(4):301-311.
908
- 909 Dos Santos Scholz, M.B., Kitzberger, C.S.G., Pagiatto, N.F., Pereira, L.F.P., Davrieux, F., Pot, D.,
910 Charmetant, P. and Leroy, T., 2016. Chemical composition in wild Ethiopian Arabica
911 coffee accessions. *Euphytica*: 429-438.
912
- 913 Dos Santos Scholz, M.B., Kitzberger, C.S.G., Pereira, L.F.P., Davrieux, F., Pot, D., Charmetant,
914 P. and Leroy, T., 2014. Application of near infrared spectroscopy for green coffee
915 biochemical phenotyping. *Journal of Near Infrared Spectroscopy*, 22(6):411-421.
916
- 917 Dotto, A.C., Dalmolin, R.S.D., ten Caten, A. and Grunwald, S., 2018. A systematic study on the
918 application of scatter-corrective and spectral-derivative preprocessing for multivariate
919 prediction of soil organic carbon by Vis-NIR spectra. *Geoderma*. 314:262-274.
920
- 921 Dubale, P., 1996. Availability of phosphorus in the coffee soil of South West Ethiopia.
922
- 923 ECX (Ethiopian Commodity Exchange), 2011. ECX quality operation manual, Addis
924 Ababa, Ethiopia.
925
- 926 Esteban-Díez, I., González-Sáiz, J.M. and Pizarro, C., 2004. Prediction of sensory properties of
927 espresso from roasted coffee samples by near-infrared spectroscopy. *Analytica Chimica*
928 *Acta*. 525(2):171-182.
929
- 930 Freeman, E.A., Moisen, G.G., Coulston, J.W. and Wilson, B.T., 2016. Random forests and
931 stochastic gradient boosting for predicting tree canopy cover: comparing tuning processes
932 and model performance. *Canadian Journal of Forest Research*, 46(3):323-339.
933
- 934 Genisheva, Z., Quintelas, C., Mesquita, D.P., Ferreira, E.C., Oliveira, J.M. and Amaral, A.L.,
935 2018. New PLS analysis approach to wine volatile compounds characterization by near
936 infrared spectroscopy (NIR). *Food chemistry*. 246:172-178.
937
- 938 Geromel, C., Ferreira, L.P., Davrieux, F., Guyot, B., Ribeyre, F., dos Santos Scholz, M.B., Pereira,
939 L.F.P., Vaast, P., Pot, D., Leroy, T. and Androcioli Filho, A. 2008. Effects of shade on the
940 development and sugar metabolism of coffee (*Coffea arabica* L.) fruits. *Plant Physiol*
941 *Bioch*:569-579.
942
- 943 Getachew, M., Verheyen, K., Tolessa, K., Ayalew, B., Hylander, K., Tack, A., Garedew, W.,
944 Bauters, M., Boeckx, P. and De Frenne, P., 2022. Shade tree canopy cover affects coffee

- 945 plant traits across elevations in coffee farms in southwest Ethiopia. *Nordic Journal of*
946 *Botany*, 2022(2):e03383.
- 947
- 948 Ghasemi, J.B. and Tavakoli, H., 2013. Application of random forest regression to spectral
949 multivariate calibration. *Analytical Methods*, 5(7):1863-1871.
- 950
- 951 Giraud, A., Grassi, S., Savorani, F., Gavoci, G., Casiraghi, E. and Geobaldo, F., 2019.
952 Determination of the geographical origin of green coffee beans using NIR spectroscopy
953 and multivariate data analysis. *Food Control*. 99:137-145.
- 954
- 955 Greenwell, B.M., Boehmke, B.C. and Gray, B., 2020. Variable Importance Plots-An Introduction
956 to the vip Package. *R J.*, 12(1):343.
- 957
- 958 Guy, F., Prache, S., Thomas, A., Bauchart, D. and Andueza, D., 2011. Prediction of lamb meat
959 fatty acid composition using near-infrared reflectance spectroscopy (NIRS). *Food*
960 *Chemistry*, 127(3):1280-1286.
- 961
- 962 Hoffman, 2021. VariancePartition: Quantifying and interpreting drivers of variation in multilevel
963 gene expression experiments.
- 964
- 965 Janitza, S. and Hornung, R., 2018. On the overestimation of random forest's out-of-bag error. *PloS*
966 *one*, 13(8):e0201904.
- 967
- 968 Jiao, Y., Li, Z., Chen, X. and Fei, S., 2020. Preprocessing methods for near-infrared spectrum
969 calibration. *Journal of Chemometrics*, 34(11):e3306.
- 970
- 971 Kapper, C., Klont, R.E., Verdonk, J.M.A.J. and Urlings, H.A.P., 2012. Prediction of pork quality
972 with near infrared spectroscopy (NIRS): 1. Feasibility and robustness of NIRS
973 measurements at laboratory scale. *Meat Science*, 91(3):294-299.
- 974
- 975 Kath, J., Byrareddy, V.M., Craparo, A., Nguyen-Huy, T., Mushtaq, S., Cao, L. and Bossolasco,
976 L., 2020. Not so robust: Robusta coffee production is highly sensitive to
977 temperature. *Global Change Biology*. 26(6):3677-3688.
- 978
- 979 Khan, M.A., Shah, M.I., Javed, M.F., Khan, M.I., Rasheed, S., El-Shorbagy, M.A., El-Zahar, E.R.
980 and Malik, M.Y., 2022. Application of random forest for modelling of surface water
981 salinity. *Ain Shams Engineering Journal*, 13(4):101635.
- 982
- 983 Khan, Z., Gul, N., Faiz, N., Gul, A., Adler, W. and Lausen, B., 2021. Optimal trees selection for
984 classification via out-of-bag assessment and sub-bagging. *IEEE Access*, 9:28591-28607.
- 985
- 986 Knox, N.M., Grunwald, S., McDowell, M.L., Bruland, G.L., Myers, D.B. and Harris, W.G., 2015.
987 Modelling soil carbon fractions with visible near-infrared (VNIR) and mid-infrared (MIR)
988 spectroscopy. *Geoderma*, 239:229-239.
- 989

- 990 Kuang, B., Tekin, Y. and Mouazen, A.M., 2015. Comparison between artificial neural network
991 and partial least squares for on-line visible and near infrared spectroscopy measurement of
992 soil organic carbon, pH and clay content. *Soil and Tillage Research*, 146:243-252.
993
- 994 Kucharska-Ambrożej, K. and Karpinska, J., 2020. The application of spectroscopic techniques in
995 combination with chemometrics for detection adulteration of some herbs and
996 spices. *Microchemical Journal*, 153:104278.
997
- 998 Kuhn, M. and Johnson, K., 2013. Regression trees and rule-based models. In *Applied predictive
999 modeling*.173-220. Springer, New York, NY.
1000
- 1001 Kutlug Sahin, E. and Colkesen, I., 2021. Performance analysis of advanced decision tree-based
1002 ensemble learning algorithms for landslide susceptibility mapping. *Geocarto
1003 International*, 36(11):1253-1275.
1004
- 1005 Legesse, A (2020) Assessment of coffee (*Coffea arabica* L.) genetic erosion and genetic resources
1006 management in Ethiopia. *Int J Agr Ext*: 223-229.
1007
- 1008 Leonel, L.E. and Philippe, V., 2007. Effects of altitude, shade, yield and fertilization on coffee
1009 quality (*Coffea arabica* L.) produced in agroforestry systems of the Northern Central Zones
1010 of Nicaragua. *J Food Sci*: 2356-2361.
1011
- 1012 Levate Macedo, L., Da Silva Araújo, C., Costa Vimercati, W., Gherardi Hein, P.R., Pimenta, C.J.
1013 and Henriques Saraiva, S., 2021. Evaluation of chemical properties of intact green coffee
1014 beans using near-infrared spectroscopy. *Journal of the Science of Food and
1015 Agriculture*, 101(8):3500-3507.
1016
- 1017 Lovatti, B.P., Nascimento, M.H., Neto, A.C., Castro, E.V. and Filgueiras, P.R., 2019. Use of
1018 Random forest in the identification of important variables. *Microchemical
1019 Journal*, 145:1129-1134.
1020
- 1021 MacCarthy, D.S., Kihara, J., Masikati, P. and Adiku, S.G., 2018. Decision support tools for site-
1022 specific fertilizer recommendations and agricultural planning in selected countries in sub-
1023 Sahara Africa. In *Improving the Profitability, Sustainability and Efficiency of Nutrients
1024 Through Site Specific Fertilizer Recommendations in West Africa Agro-Ecosystems*: 265-
1025 289.
1026
- 1027 Manthou, E., Lago, S.L., Dagres, E., Lianou, A., Tsakanikas, P., Panagou, E.Z., Anastasiadi, M.,
1028 Mohareb, F. and Nychas, G.J.E., 2020. Application of spectroscopic and multispectral
1029 imaging technologies on the assessment of ready-to-eat pineapple quality: A performance
1030 evaluation study of machine learning models generated from two commercial data
1031 analytics tools. *Computers and Electronics in Agriculture*, 175:105529.
1032

- 1033 Martínez-Santos, P., Aristizábal, H.F., Díaz-Alcaide, S. and Gómez-Escalonilla, V., 2021.
1034 Predictive mapping of aquatic ecosystems by means of support vector machines and
1035 random forests. *Journal of Hydrology*, 595:126026.
1036
- 1037 Mishra, P. and Woltering, E.J., 2021. Identifying key wavenumbers that improve prediction of
1038 amylose in rice samples utilizing advanced wavenumber selection
1039 techniques. *Talanta*. 224:121908.
1040
- 1041 Monteiro, L.R., Lange, C.N., Freire, B.M., Pedron, T., Da Silva, J.J.C., De Magalhães Junior,
1042 A.M., Pegoraro, C., Busanello, C. and Batista, B.L., 2020. Inter-and intra-variability in the
1043 mineral content of rice varieties grown in various microclimatic regions of southern
1044 Brazil. *Journal of Food Composition and Analysis*, 92:103535.
1045
- 1046 Muleta, D., Assefa, F., Nemomissa, S. and Granhall, U., 2008. Distribution of arbuscular
1047 mycorrhizal fungi spores in soils of smallholder agroforestry and monocultural coffee
1048 systems in southwestern Ethiopia. *Biology and fertility of soils*. 44(4):653-659.
1049
- 1050 Naccarato, A., Furia, E., Sindona, G. and Tagarelli, A., 2016. Multivariate class modeling
1051 techniques applied to multielement analysis for the verification of the geographical origin
1052 of chili pepper. *Food chemistry*, 206:217-222.
1053
- 1054 Nakagawa, S., Johnson, P.C. and Schielzeth, H., 2017. The coefficient of determination R^2 and
1055 intra-class correlation coefficient from generalized linear mixed-effects models revisited
1056 and expanded. *Journal of the Royal Society Interface*, 14(134):207-213.
1057
- 1058 Navarro-Serrano, F., López-Moreno, J.I., Azorin-Molina, C., Alonso-González, E., Aznarez-
1059 Balta, M., Buisán, S.T. and Revuelto, J., 2020. Elevation Effects on Air Temperature in a
1060 Topographically Complex Mountain Valley in the Spanish
1061 Pyrenees. *Atmosphere*, 11(6):656.
1062
- 1063 Nawar, S. and Mouazen, A.M., 2017. Comparison between random forests, artificial neural
1064 networks and gradient boosted machines methods of on-line Vis-NIR spectroscopy
1065 measurements of soil total nitrogen and total carbon. *Sensors*, 17(10):2428.
1066
- 1067 Nesper, M., Kueffer, C., Krishnan, S., Kushalappa, C.G. and Ghazoul, J., 2017. Shade tree
1068 diversity enhances coffee production and quality in agroforestry systems in the Western
1069 Ghats. *Agriculture, Ecosystems & Environment*, 247:172-181.
1070
- 1071 Ng, W., Minasny, B., Malone, B. and Filippi, P., 2018. In search of an optimum sampling
1072 algorithm for prediction of soil properties from infrared spectra. *PeerJ*, 6:e5722.
1073
- 1074 NS Akbar, M., Rachmawati, E. and Sthevanie, F., 2020. Visual Feature and Machine Learning
1075 Approach for Arabica Green Coffee Beans Grade Determination. In 2020 the 6th
1076 International Conference on Communication and Information Processing. 97-104.
1077

- 1078 P. Probst. measures: Performance Measures for Statistical Learning, 2018. URL [https://CRAN.R-](https://CRAN.R-project.org/package=measures)
1079 [project.org/package=measures](https://CRAN.R-project.org/package=measures). R package version 0.2.
1080
- 1081 P. Probst. varImp: RF Variable Importance for Arbitrary Measures, 2019. URL [https://CRAN.R-](https://CRAN.R-project.org/package=varImp)
1082 [project.org/package=varImp](https://CRAN.R-project.org/package=varImp). R package version 0.3.
1083
- 1084 Okubo, N. and Kurata, Y., 2019. Nondestructive classification analysis of green coffee beans by
1085 using near-infrared spectroscopy. *Foods*, 8(2):82.
1086
- 1087 Onwuka, B. and Mang, B., 2018. Effects of soil temperature on some soil properties and plant
1088 growth. *Adv Plants Agric Res*. 8(1):34.
1089
- 1090 Ovalle-Rivera, O., Laderach, P., Bunn, C., Obersteiner, M. and Schroth, G., 2015. Projected shifts
1091 in *Coffea arabica* suitability among major global producing regions due to climate change.
1092 *Plos One*: 24.
1093
- 1094 Oyinbo, O., Chamberlin, J., Vanlauwe, B., Vranken, L., Kamara, Y.A., Craufurd, P. and Maertens,
1095 M., 2019. Farmers' preferences for high-input agriculture supported by site-specific
1096 extension services: Evidence from a choice experiment in Nigeria. *Agricultural*
1097 *systems*, 173:12-26.
1098
- 1099 Peng, X., Shi, T., Song, A., Chen, Y. and Gao, W., 2014. Estimating soil organic carbon using
1100 VIS/NIR spectroscopy with SVMR and SPA methods. *Remote Sensing*, 6(4):2699-2717.
1101
- 1102 Pizarro, C., Esteban-Díez, I., González-Sáiz, J.M. and Forina, M., 2007. Use of near-infrared
1103 spectroscopy and feature selection techniques for predicting the caffeine content and
1104 roasting color in roasted coffees. *Journal of Agricultural and Food*
1105 *Chemistry*, 55(18):7477-7488.
1106
- 1107 Probst, P., Wright, M.N. and Boulesteix, A.L., 2019. Hyperparameters and tuning strategies for
1108 random forest. *Wiley Interdisciplinary Reviews: data mining and knowledge*
1109 *discovery*, 9(3):e1301.
1110
- 1111 R Core Team (2022). R: A language and environment for statistical computing. R Foundation for
1112 Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.
1113
- 1114 Ren, G., Ning, J. and Zhang, Z., 2021. Multi-variable selection strategy based on near-infrared
1115 spectra for the rapid description of dianhong black tea quality. *Spectrochimica Acta Part*
1116 *A: Molecular and Biomolecular Spectroscopy*, 245:118918.
1117
- 1118 Revuelto, J., Azorin-Molina, C., Alonso-González, E., Sanmiguel-Valladolid, A., Navarro-Serrano,
1119 F., Rico, I. and López-Moreno, J.I., 2017. Meteorological and snow distribution data in the
1120 Izas Experimental Catchment (Spanish Pyrenees) from 2011 to 2017. *Earth System*
1121 *Science Data*, 9(2):993-1005.
1122

- 1123 Ribeiro, J.S., Ferreira, M.M. and Salva, T.J.G., 2011. Chemometric models for the quantitative
1124 descriptive sensory analysis of Arabica coffee beverages using near infrared
1125 spectroscopy. *Talanta*. 83(5):1352-1358.
1126
- 1127 Ribeiro, J.S., Salva, T.D.J.G. and Silvarolla, M.B., 2021. Prediction of a wide range of compounds
1128 concentration in raw coffee beans using NIRS, PLS and variable selection. *Food*
1129 *Control*, 125:107967.
1130
- 1131 Robertson, G.P., Coleman, D.C., Sollins, P. and Bledsoe, C.S. eds., 1999. *Standard soil methods*
1132 *for long-term ecological research*. Oxford University Press on Demand.
1133
- 1134 Rodríguez-Pulido, F.J., Barbin, D.F., Sun, D.W., Gordillo, B., González-Miret, M.L. and Heredia,
1135 F.J., 2013. Grape seed characterization by NIR hyperspectral imaging. *Postharvest Biology*
1136 *and Technology*, 76:74-82.
1137
- 1138 Sarmiento-Soler, A., Rötter, R.P., Hoffmann, M.P., Jassogne, L., van Asten, P., Graefe, S. and
1139 Vaast, P., 2022. Disentangling effects of altitude and shade cover on coffee fruit dynamics
1140 and vegetative growth in smallholder coffee systems. *Agriculture, Ecosystems &*
1141 *Environment*, 326:107786.
1142
- 1143 Semedo, J.N., Rodrigues, W.P., Dubberstein, D., Martins, M.Q., Martins, L.D., Pais, I.P.,
1144 Rodrigues, A.P., Leitão, A.E., Partelli, F.L., Campostrini, E. and Tomaz, M.A., 2018.
1145 Coffee responses to drought, warming and high [CO₂] in a context of future climate
1146 change scenarios. In *Theory and practice of climate adaptation*. 465-477.
1147
- 1148 Sida, T.S., Chamberlin, J., Ayalew, H., Kosmowski, F. and Craufurd, P., 2021. Implications of
1149 intra-plot heterogeneity for yield estimation accuracy: Evidence from smallholder maize
1150 systems in Ethiopia. *Field crops research*, 267:108147.
1151
- 1152 Silva, E.A.D., Mazzafera, P., Brunini, O., Sakai, E., Arruda, F.B., Mattoso, L.H.C., Carvalho, C.R.
1153 and Pires, R.C.M., 2005. The influence of water management and environmental
1154 conditions on the chemical composition and beverage quality of coffee beans. *Plant*
1155 *physiol.* 229-238.
1156
- 1157 Silva, T.V., Hubinger, S.Z., Neto, J.A.G., Milori, D.M.B.P., Ferreira, E.J. and Ferreira, E.C., 2017.
1158 Potential of Laser Induced Breakdown Spectroscopy for analyzing the quality of unroasted
1159 and ground coffee. *Spectrochimica Acta Part B: Atomic Spectroscopy*, 135:29-33.
1160
- 1161 Somporn, C., Kantuo, A., Theerakulpisut, P. and Siriamornpun, S., 2012. Effect of shading on
1162 yield, sugar content, phenolic acids and antioxidant property of coffee beans (*Coffea*
1163 *arabica* L.) harvested from north-eastern Thailand. *Journal of the Science of Food and*
1164 *Agriculture*, 92(9):1956-1963.
1165

- 1166 Souza, J.C., Pasquini, C. and Hespanhol, M.C., 2022. Feasibility of compact near-infrared
1167 spectrophotometers and multivariate data analysis to assess roasted ground coffee
1168 traits. *Food Control*. 109041.
1169
- 1170 Teklu, B., Mohammed, A. and Kufa, T., 2011. Effect of processing methods and drying materials
1171 on bean physical and sensorial quality attributes of coffee (*Coffea arabica* L.) Varieties at
1172 Gera and Jimma (Doctoral dissertation).
1173
- 1174 Teles, G., Rodrigues, J.J., Rabelo, R.A. and Kozlov, S.A., 2021. Comparative study of support
1175 vector machines and random forests machine learning algorithms on credit
1176 operation. *Software: Practice and Experience*, 51(12):2492-2500.
1177
- 1178 Tolessa, K., Dheer, J., Duchateau, L. and Boeckx, P., 2017. Influence of growing altitude, shade
1179 and harvest period on quality and biochemical composition of Ethiopian specialty coffee.
1180 *J Sci Food Agr*. 2849-2857.
1181
- 1182 Tolessa, K., Rademaker, M., De Baets, B. and Boeckx, P., 2016. Prediction of specialty coffee cup
1183 quality based on near infrared spectra of green coffee beans. *Talanta*, 150:367-374.
1184
- 1185 Trevisan, R.G., Bullock, D.S. and Martin, N.F., 2021. Spatial variability of crop responses to
1186 agronomic inputs in on-farm precision experimentation. *Precision Agriculture*, 22(2):342-
1187 363.
1188
- 1189 Tridawati, A., Wikantika, K., Susantoro, T.M., Harto, A.B., Darmawan, S., Yayusman, L.F. and
1190 Ghazali, M.F., 2020. Mapping the distribution of coffee plantations from multi-resolution,
1191 multi-temporal, and multi-sensor data using a random forest algorithm. *Remote*
1192 *Sensing*, 12(23):3933.
1193
- 1194 Tyrallis, H., Papacharalampous, G. and Langousis, A., 2019. A brief review of random forests for
1195 water scientists and practitioners and their recent history in water resources. 11(5):910.
1196
- 1197 Van Loon, M.P., Adjei-Nsiah, S., Descheemaeker, K., Akotsen-Mensah, C., van Dijk, M., Morley,
1198 T., van Ittersum, M.K. and Reidsma, P., 2019. Can yield variability be explained?
1199 Integrated assessment of maize yield gaps across smallholders in Ghana. *Field Crops*
1200 *Research*. 236:132-144.
1201
- 1202 Vargas, C. and El Hanandeh, A., 2021. Systematic literature review, meta-analysis and artificial
1203 neural network modelling of plastic waste addition to bitumen. *Journal of Cleaner*
1204 *Production*, 280:124369.
1205
- 1206 Viscarra Rossel, R.A. and Lark, R.M., 2009. Improved analysis and modelling of soil diffuse
1207 reflectance spectra using wavelets. *European Journal of Soil Science*, 60(3):453-464.
1208
- 1209 Wadoux, A.M.C., Brus, D.J. and Heuvelink, G.B., 2019. Sampling design optimization for soil
1210 mapping with random forest. *Geoderma*, 355:113913.

- 1211
1212 Wadoux, A.M.C., Malone, B., Minasny, B., Fajardo, M. and McBratney, A.B., 2021. Soil Spectral
1213 Inference with R: Analyzing Digital Soil Spectra Using the R Programming Environment.
1214 Springer Nature.
1215
- 1216 Wang, L., Li, Q., Yu, Y. and Liu, J., 2018. Region compatibility-based stability assessment for
1217 decision trees. *Expert Systems with Applications*, 105:112-128.
1218
- 1219 Worku, M., De Meulenaer, B., Duchateau, L. and Boeckx, P., 2018. Effect of altitude on
1220 biochemical composition and quality of green arabica coffee beans can be affected by
1221 shade and postharvest processing method. *Food Res Int*, 105:278-285.
1222
- 1223 Wright, M.N., Dankowski, T. and Ziegler, A., 2017. Unbiased split variable selection for random
1224 survival forests using maximally selected rank statistics. *Statistics in*
1225 *medicine*, 36(8):1272-1284.
1226
- 1227 Xie, X., Wu, T., Zhu, M., Jiang, G., Xu, Y., Wang, X. and Pu, L., 2021. Comparison of random
1228 forest and multiple linear regression models for estimation of soil extracellular enzyme
1229 activities in agricultural reclaimed coastal saline land. *Ecological Indicators*. 120:106925.
1230
- 1231 Yadessa, A., Burkhardt, J., Bekele, E., Hundera, K. and Goldbach, H., 2020. Influence of soil
1232 properties on bean quality of wild *Coffea arabica* in the natural coffee forests of southwest
1233 and southeast Ethiopia. *Ethiopian J App Sci Tech*. 23-38.
1234
- 1235 Yergenson, N. and Aston, D.E., 2020. Monitoring coffee roasting cracks and predicting with in
1236 situ near-infrared spectroscopy. *Journal of Food Process Engineering*, 43(2):e13305.
1237
- 1238 Zhang, C., Jiang, H., Liu, F. and He, Y., 2017. Application of near-infrared hyperspectral imaging
1239 with variable selection methods to determine and visualize caffeine content of coffee
1240 beans. *Food and bioprocess technology*, 10(1):213-221.
1241
- 1242 Zhou, H., Zhang, J., Zhou, Y., Guo, X. and Ma, Y., 2021. A feature selection algorithm of decision
1243 tree based on feature weight. *Expert Systems with Applications*, 164:113842.
1244
- 1245 Zhu, M., Long, Y., Chen, Y., Huang, Y., Tang, L., Gan, B., Yu, Q. and Xie, J., 2021. Fast
1246 determination of lipid and protein content in green coffee beans from different origins using
1247 NIR spectroscopy and chemometrics. *Journal of Food Composition and*
1248 *Analysis*, 102:104055.