# Correlation network analysis based on untargeted LC-MS profiles of cocoa reveals processing stage and origin country — Source link ↗

Santhust Kumar, Roy N. Dsouza, Marcello Corno, Matthias S. Ullrich ...+2 more authors

**Institutions:** Jacobs University Bremen

# Correlation network analysis based on untargeted LC-MS profiles of cocoa reveals processing stage and origin country

Santhust Kumar,[1*] Roy N. D'Souza,[1] Marcello Corno,[2] Matthias S. Ullrich,[1] Nikolai Kuhnert[1] and Marc-Thorsten Hütt[1*]

[1]Department of Life Sciences and Chemistry, Jacobs University Bremen, Campus Ring 1, 28759 Bremen, Germany

[2]Barry Callebaut AG, Westpark, Pfingstweidstrasse 60, Zurich 8005, Switzerland

*Correspondence to:

Dr. Santhust Kumar (s.santhust@jacobs-university.de)

Prof. Dr. Marc-Thorsten Hütt (m.huett@jacobs-university.de)

Department of Life Sciences & Chemistry

Jacobs University Bremen gGmbH

Campus Ring 1

28759 Bremen, Germany

1

20    ABSTRACT

21    In order to implement quality control measures and create fine flavor products, an important

22    objective in cocoa processing industry is to realize standards for characterization of cocoa raw

23    materials, intermediate and finished products with respect to their processing stages and

24    countries of origin. Towards this end, various works have studied separability or

25    distinguishability of cocoa samples belonging to various processing stages in a typical cocoa

26    processing pipeline or to different origins. Limited amount of success has been possible in this

27    direction in that unfermented and fermented cocoa samples have been shown to group into

28    separate clusters in PCA. However, a clear clustering with respect to the country of origin has

29    remained elusive. In this work we suggest an alternative approach to this problem through the

30    framework of correlation networks. For 140 cocoa samples belonging to eight countries and

31    three progressive stages in a typical cocoa processing pipeline we compute pairwise Spearman

32    and Pearson correlation coefficients based on the LC-MS profiles and derive correlation

33    networks by retaining only correlations higher than a threshold. Progressively increasing this

34    threshold reveals, first, processing stage (or sample type) modules (or network clusters) at low

35    and intermediate values of correlation threshold and then country specific modules at high

36    correlation thresholds. We present both qualitative and quantitative evidence through network

37    visualization and node connectivity statistics. Besides demonstrating separability of the two

38    data properties via this network-based method, our work suggests a new approach for studying

39    classification of cocoa samples with nested attributes of processing stage sample types and

40    country of origin along with possibility of including additional factors, e.g., hybrid variety, etc.

41    in the analysis.

42

2

43    **Keywords:** *Theobroma cacao*, LC-MS, correlation network, origin classification, processing-

44    stage classification.

45

3

## 1. Introduction

Cocoa, scientifically *Theobroma cacao*, is a commodity of commercial interest to farmers as a crop and to businesses as a raw material for producing various cocoa based food products. Therefore, quality, variety and characteristics of cocoa and its derived food items have become an important area of research and development. Quality control (Fayeulle et al., 2019; Guehi et al., 2010; Kongor et al., 2016; Lima et al., 2011) and design of single origin cocoa products (Oberrauter et al., 2018; Ozretic-Dosen et al., 2007) are two of many focus areas in cocoa research. The former helps in ensuring whether the stages in a typical cocoa processing pipeline have been rightly carried to achieve the best possible finished product, and the latter commands high value among consumers for nuanced taste and aroma of the consumed food item.

Previous research successfully demonstrate characteristic differences between unfermented, partially fermented and fermented cocoa samples (processing-stages) and even identified corresponding potentially responsible classes of compounds through multivariate statistical analysis, e.g., principal component analysis (PCA) (Wold et al., 1987), on the chemical composition of these samples (Caligiani et al., 2014; D'Souza et al., 2017; Kumari et al., 2018; Megías-Pérez et al., 2018). Baring a few cases where the number of distinct countries relating to the samples in dataset at hand is few (D'Souza et al., 2017; Milev et al., 2014; Oliveira et al., 2016) or based on large continental regions (Acierno et al., 2016, 2018; Bertoldi et al., 2016; Kumari et al., 2018; Marseglia et al., 2016), a successful characteristic differentiation amongst samples on the basis of their country of origin has remained hard to define through metabolomic analysis (D'Souza et al., 2017; Sirbu et al., 2018; Vázquez-Ovando et al., 2015).

On the other hand, the 'language of networks' (Albert and Barabási, 2002; Newman, 2003) has proven immensely useful in visualizing and interpreting relationships between multitude of entities, and across many disciplines—metabolomics (Jeong et al., 2000), genetics (Grimbs et

4

70    al., 2019; Kumar et al., 2018), proteomics (Szklarczyk et al., 2015), social science (Borgatti et

71    al., 2009), logistics (Becker et al., 2012), gut ecology (Claussen et al., 2017) medicine

72    (Barabási et al., 2011; Batushansky et al., 2016), finance (Kumar and Deo, 2012; Namaki et

73    al., 2011), etc. to name a few. Some works have successfully applied this approach in the field

74    of food science (Ahn et al., 2011; Hochberg et al., 2013; Ursem et al., 2008; Wang et al., 2017).

75    Here, we apply the framework of network science to simultaneously study the clustering of

76    cocoa samples with regards to their processing-stage sample types and country of origin.

77    We start by computing pairwise Spearman and Pearson correlation coefficients between 140

78    cocoa samples belonging to three different stages in a typical cocoa processing pipeline

79    (unfermented, fermented and liquor) and 8 countries through their LC-MS profiles in positive

80    ion mode. On the basis of correlations obtained, we construct correlation networks, at varying

81    correlation thresholds. In these networks, the nodes are samples and an edge between two

82    samples is drawn, when the correlation coefficient exceeds the threshold value.

83    We find that, as we progressively increase the correlation threshold from 0 towards 1, the

84    clustering of cocoa samples is first dominated by processing-stage sample types at low and

85    intermediate correlation thresholds, and then by countries of origin at high correlation

86    thresholds. We show this both qualitatively and quantitatively via network visualizations and

87    network edge statistics.

88    Our work demonstrates the presence of processing-stage level grouping on a coarser level and

89    origin level grouping on a finer level within the former. This nested grouping can be revealed

90    by successively keeping higher correlations. Further, our works suggests a new approach to

91    study clustering or classification of food samples upon multiple nested attributes and can prove

92    an important complement to traditional approaches and strategies.

93

5

## 2. Materials and Methods

### 2.1 Country and Origin details

The LC-MS data set we use here has a total of 140 samples (positive ion mode). The samples have been gathered and their LC-MS profiling done under COMETA project over a range of about past five years. These samples can be grouped into three sample-types (Unfermented, Fermented and Liquors) and eight origins (Brazil, Cameroon, Ecuador, Ghana, Indonesia, Ivory Coast, Malaysia and Tanzania). A cross-table of details about number of samples belonging to particular sample-type and country is given in Table 1.

| | Brazil | Cameroon | Ecuador | Ghana | Indonesia | Ivory coast | Malaysia | Tanzania | All |
|---|---|---|---|---|---|---|---|---|---|
| Unfermented | 4 | 3 | 8 | 0 | 14 | 16 | 6 | 3 | 54 |
| Fermented | 4 | 3 | 12 | 0 | 16 | 16 | 3 | 9 | 63 |
| Liquor | 0 | 6 | 3 | 5 | 0 | 9 | 0 | 0 | 23 |
| All | 8 | 12 | 23 | 5 | 30 | 41 | 9 | 12 | 140 |

**Table 1 Sample division.** The LCMS data set can be grouped on twin axes: sample-type and origin. There are 3 sample-types: Unfermented, Fermented and Liquors, and there are 8 origins (Brazil, Cameroon, Ecuador, Ghana, Indonesia, Ivory Coast, Malaysia and Tanzania).

### 2.2 Data pre-processing and cleaning

The data generation, cleaning, standardization and organization has been discussed in an earlier work (Kumar et al *previous manuscript*). Briefly, LC-MS data of all the samples was processed

6

110  using MZMine (Pluskal et al., 2010) giving peak area list and corresponding *m/z ratio* and

111  *retention times*. The detected compounds are assigned names/chemical formula on the basis of

112  four ionization states ([M+H], [M+2H], [M+3H], [2M+H]) when possible, else the compound

113  was named as 'Unknown_' suffixed with the *m/z* value, e.g., Unknown_865.1927. The samples

114  were then put in an excel file, where each row represents a sample, and the column contain

115  information about the sample-type, origin and peak areas of various compounds sorted in

116  descending order by their mean peak are across all the samples.

## 2.3 Network production and visualization

118  Spearman and Pearson correlation analysis, and network generation/transformation was carried

119  by writing programs from scratch in Python programming language making use of popular

120  modules such as Pandas (McKinney, 2010, 2011) and NetworkX (Hagberg et al., 2008).

121  Network visualization has been done in Cytoscape (Shannon et al., 2003). For layout of the

122  network either of the following two variants of spring layout, which were available in

123  Cytoscape itself, were used: (a) Edge-weighted Spring Embedded Layout (Kamada and Kawai,

124  1989), (b) Compound Spring Embedder (CoSE) (Dogrusoz et al., 2009). These layouts take

125  into account the weight of the edge (in our case the Spearman or Pearson correlation

126  coefficient) between nodes, so that the nodes with higher weight (correlations) are placed closer

127  together.

## 2.4 Null model network or control network

129  A null model network is made by randomizing the weights (correlations) of edges in the

130  original correlation network. It is important to note that the null model network so obtained has

131  the same correlation distribution as that of the original correlation network because the set of

132  correlations in the network remains unchanged, only the correlations between nodes is

133  randomized. An ensemble of 100 such null model networks were generated. The reported

7

134    statistics about a studied property on the null model networks is obtained by making

135    calculations over this ensemble and then reporting the mean and standard deviation of the

136    studied property. Higher the difference in the studied property between the original network

137    and null network ensemble, higher the significance of the observed property in the original

138    network.

## 3. Results

### 3.1 Correlation between cocoa samples

141    A typical LC-MS profile contains information about thousands of compounds present in a

142    given sample defined by their retention time and associated *m/z* values (Kuhnert et al., 2013).

143    Using the areas of peaks as a rough measure for concentration of these compounds across all

144    samples, we calculate the Spearman and Pearson correlation coefficients (*r*) for all pairs of

145    samples in our dataset.

146    The LC-MS data can be represented as a matrix L with entries $l_i^\alpha$. The upper index $\alpha$ represents

147    the sample and lower index *i* represents the compound. Thus, the scalar quantity $l_i^\alpha$ represents

148    the concentration of $i^{\text{th}}$ compound in the $\alpha^{\text{th}}$ LC-MS sample. Correspondingly, $l^\alpha$ is a vector

149    which represents the LC-MS profile of sample $\alpha$. The Pearson correlation between two LC-

150    MS samples, say $\alpha$ and $\beta$ with corresponding profiles $l^\alpha$ and $l^\beta$, can be denoted as $r_{\alpha\beta}$. It is

151    calculated as

152
$$r_{\alpha\beta} = \frac{\text{cov}(l^\alpha, l^\beta)}{\sigma_{l^\alpha}\sigma_{l^\beta}}$$
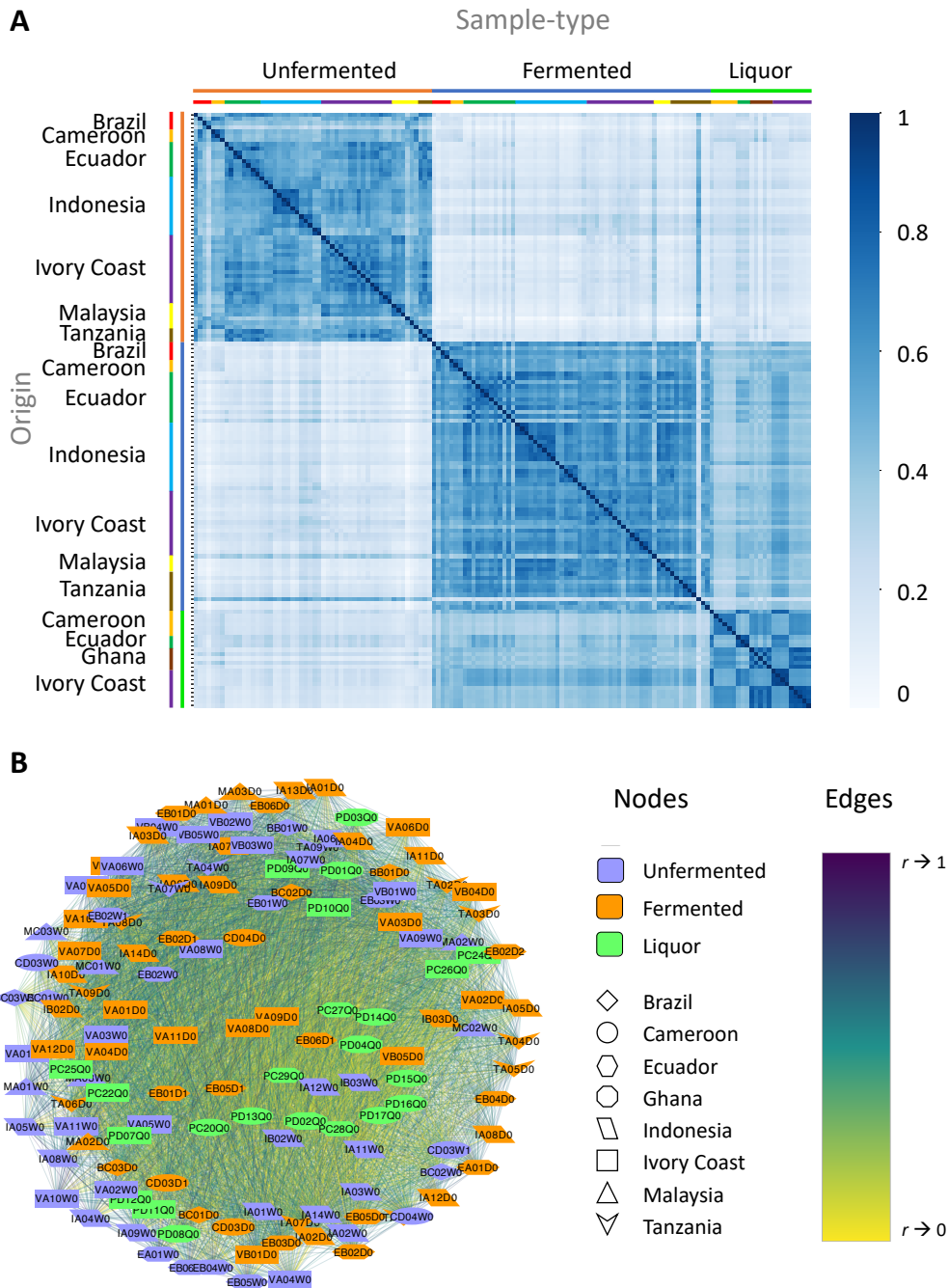
153    Where $\text{cov}(l^\alpha, l^\beta)$ represents the covariance between the LC-MS profiles of samples $\alpha$ and $\beta$,

154    while $\sigma_{l^\alpha}$ and $\sigma_{l^\beta}$ represent the standard deviation in the LC-MS profiles $l^\alpha$ and $l^\beta$,

155    respectively. The Spearman correlation can be defined as the Pearson correlation between the

8

156     ranks of the original variables (i.e., $l^\alpha$ and $l^\beta$). The ranked variables $\tilde{l}^\alpha$ and $\tilde{l}^\beta$, are obtained

157     from the original variables $l^\alpha$ and $l^\beta$ by sorting them from lowest to highest and substituting

158     the values by the position in the sorted list (i.e., the rank of the values). Formally, the Spearman

159     correlation coefficient is thus calculated as

160    
$$\tilde{r}_{\alpha\beta} = \frac{\mathrm{cov}(\tilde{l}^\alpha, \tilde{l}^\beta)}{\sigma_{\tilde{l}^\alpha} \sigma_{\tilde{l}^\beta}}$$

161     The Spearman and Pearson correlations across all pairs of LC-MS samples can be written in

162     the form of matrices, $\widetilde{R}$ and $R$, whose entries denoted by $\tilde{r}_{\alpha\beta}$ and $r_{\alpha\beta}$, respectively.

163     The correlation matrices $\widetilde{R}$ so obtained, i.e. the case of Spearman correlation coefficient, is

164     visualized through heatmap in Figure 1A. The heat map of Pearson correlation coefficient

165     matrix, $R$, is given in Supplementary Information file. By construction the correlation matrices

166     $\widetilde{R}$ and $R$ are symmetric. The twin attributes of nodes, namely the processing-stage sample type

167     and country of origin, have been alternatively marked on the sides. Three blocks corresponding

168     to Unfermented, Fermented and Liquor samples blocks are clearly distinguishable. It is also

169     visible that Fermented and Liquor samples are part of a larger block which is separated from

170     Unfermented samples. This shows that Liquor samples are closer in character to Fermented

171     samples. This is in consonance with general expectation that liquor follows the fermentation

172     stage. Furthermore, more chemical changes occur in cocoa when moving from unfermented

173     stage to fermented stage than occurs from fermented to liquor stage. In case of correlation

174     heatmap obtained using Pearson correlation (Supplementary Information file) the block of

175     Unfermented samples is clearly distinguishable from Fermented and Liquor samples, while the

176     Fermented and Liquor samples are mildly distinguishable. Further, it is important to note that

177     no block structure on the basis of country is discernable at this level of detail about the

178     correlations.

      9

**Figure 1 Correlation between cocoa samples. (A)** Correlation heatmap**.** Darker regions represent high correlation, and lighter regions represent low correlation. Samples have been sorted on twin axes, first on processing stage sample-[...]on country of origin. Two distinct square block regions are clearly visible along the diagonal of the matrix, corresponding to Unfermented (smaller block) and Fermented (bigger block) samples. **(B)** Correlation Network. The correlation network made using all correlations between the set of

10

186    cocoa samples using Spearman correlation. The nodes are color coded according to their

187    processing-stage sample type and shape coded by their country of origin. The colors of edges

188    code for the strength of correlation between nodes. The network is visualized using Cytoscape

189    (Shannon et al., 2003) with 'edge-weighted spring embedded layout' which keeps nodes

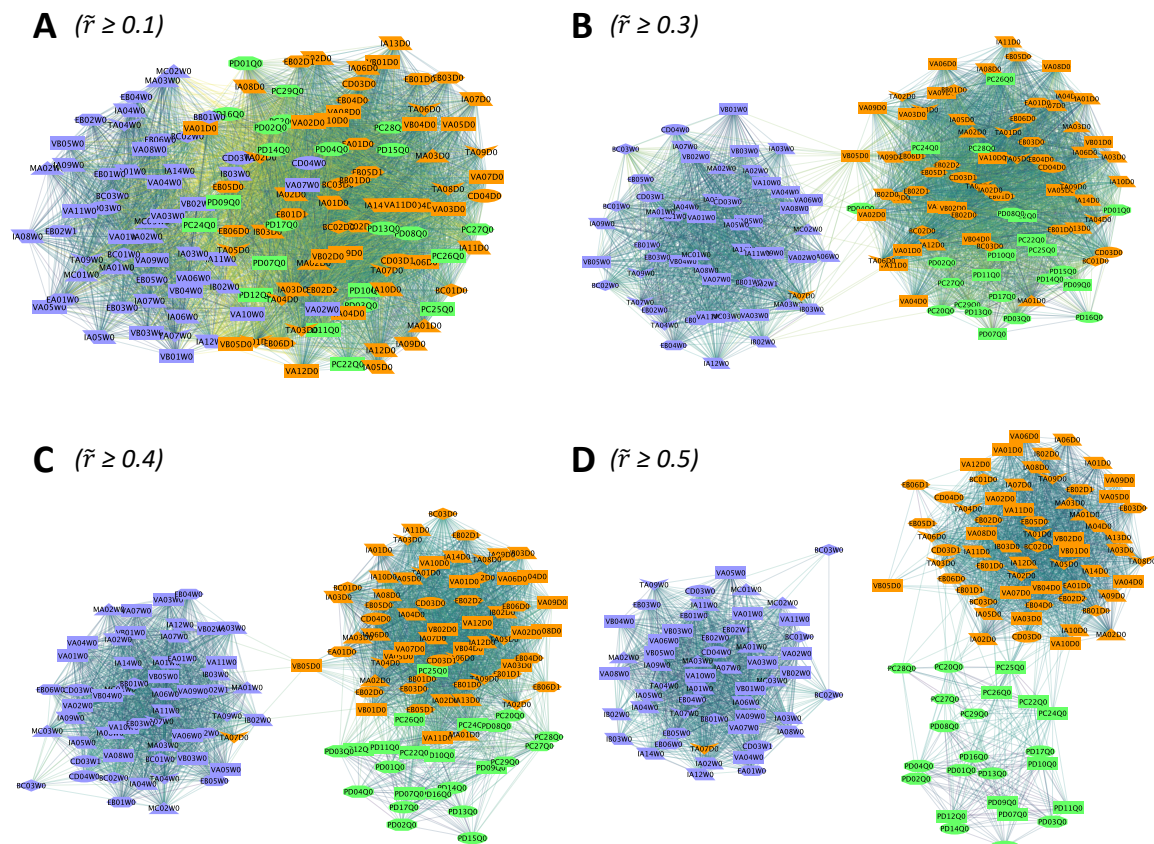190    connected with higher correlations closer together.

191

192    Next, we define correlation network using the Spearman ($\tilde{r}$) and Pearson correlations ($r$)

193    obtained above. A network is defined through two sets of entities: nodes ($N$) and edges ($E$).

194    The nodes denote the objects which are related to each other in some way, and the edge

195    represent the relation between the nodes. For further knowledge about network, see (Albert and

196    Barabási, 2002; Newman, 2003). In a correlation network, an edge represents the correlation

197    between two nodes. In our correlation network, the nodes represent the different LC-MS

198    samples of cocoa or its products sourced from different origins, and the edge between the nodes

199    represent the correlation between the LC-MS samples. Figure 1B shows the correlation

200    network obtained by using all correlations (0 to 1) between all LC-MS samples and visualized

201    with edge-weighted spring layout (see section 2.3 Network production and visualization).

202    Metadata about the LC-MS samples, such as country, and processing-stage sample type

203    (unfermented, fermented, or liquor) has been represented through color and shape of nodes,

204    respectively. The network shown in Figure 1B is the correlation network made using Spearman

205    correlation and has 140 nodes and 6833 edges, i.e. 140 cocoa LC-MS samples and 6833

206    correlations ($\tilde{r} > 0$) between the nodes. The network made using Pearson correlation is shown

207    in the Supplementary Information file. The label of the node represents the internal LC-MS id.

208    The strength of correlation is represented by the color of the edge between the nodes, yellow

209    representing low correlation and violet representing high correlation. The spatial placement of

210    nodes in Figure 1B, and all of the following networks, is done through variants of spring layout

11

211    algorithms in Cytoscape (Shannon et al., 2003) which places the nodes with higher correlation

212    closer together (2.3 Network production and visualization).

213    **3.2 Networks at low and intermediate correlation thresholds reveals processing-stage**

214    **sample type modules**

215    Next, we analyze correlation networks at low and intermediate correlation thresholds ($\tilde{r}_{th}$),

216    varying it from $\tilde{r}_{th} = 0.1$ to $0.5$, in steps of $0.1$. The network at a given correlation threshold

217    contains all the edges with correlation greater than or equal to the set threshold. Some of these

218    networks are visualized in Figure 2. Panels A, B, C and D in Figure 2 show the network at

219    correlation thresholds of 0.1, 0.3, 0.4 and 0.5, respectively. In panel A, the nodes belonging to

220    Unfermented samples are seen little separated from the nodes belonging to the Fermented and

221    Liquor samples. In panel B, the Unfermented samples are clearly separated from the Fermented

222    and Liquor samples. Within the Fermented and Liquor samples little grouping starts to form.

223    In panel C, the separation between the Fermented and Liquor samples becomes enough clear.

224    And in panel D, all the three samples can be seen clearly separated from one another. This

225    separation of samples first into two groups: (a) Unfermented, and (b) Fermented and Liquor

226    samples, and then slowly into three groups: Unfermented, Fermented and Liquor samples, is

227    in congruence with the earlier result seen in the structure of the correlation matrix heatmap

228    shown in Figure 1A. Both Figure 1A and Figure 2B,C show that the liquor sample are more

229    similar to the fermented samples than to the unfermented samples. This is in accordance with

230    the fact that major chemical and physical changes in cocoa beans takes place during the

231    processes of fermentation. A movie of the network as a function of progressively increasing

232    the threshold is attached as supplementary information which clearly shows the evolving

233    network and separation of samples belonging to different cocoa processing stage. Similar

234    behavior is noted for the case of correlation network formed using the Pearson correlation

12

235    coefficient (Supplementary Information file) however at different values of correlation

236    threshold.



**Figure 2 Processing-stage modules in low and intermediate threshold correlation networks.** The figure reveals modules of samples belonging to the same cocoa processing-stage in a typical cocoa processing pipeline. **(A)** Network of LC-MS samples at a correlation threshold of 0.1 revealing separation of unfermented, fermented and liquor cluster. **(B, C)** Correlation thresholds 0.3 and 0.4. The separation between different processing-stage sample types improves. **(D)** Correlation threshold 0.5. Three groups of unfermented, fermented and liquor samples are clearly separated. The figure follows same legend as of Figure 1B**E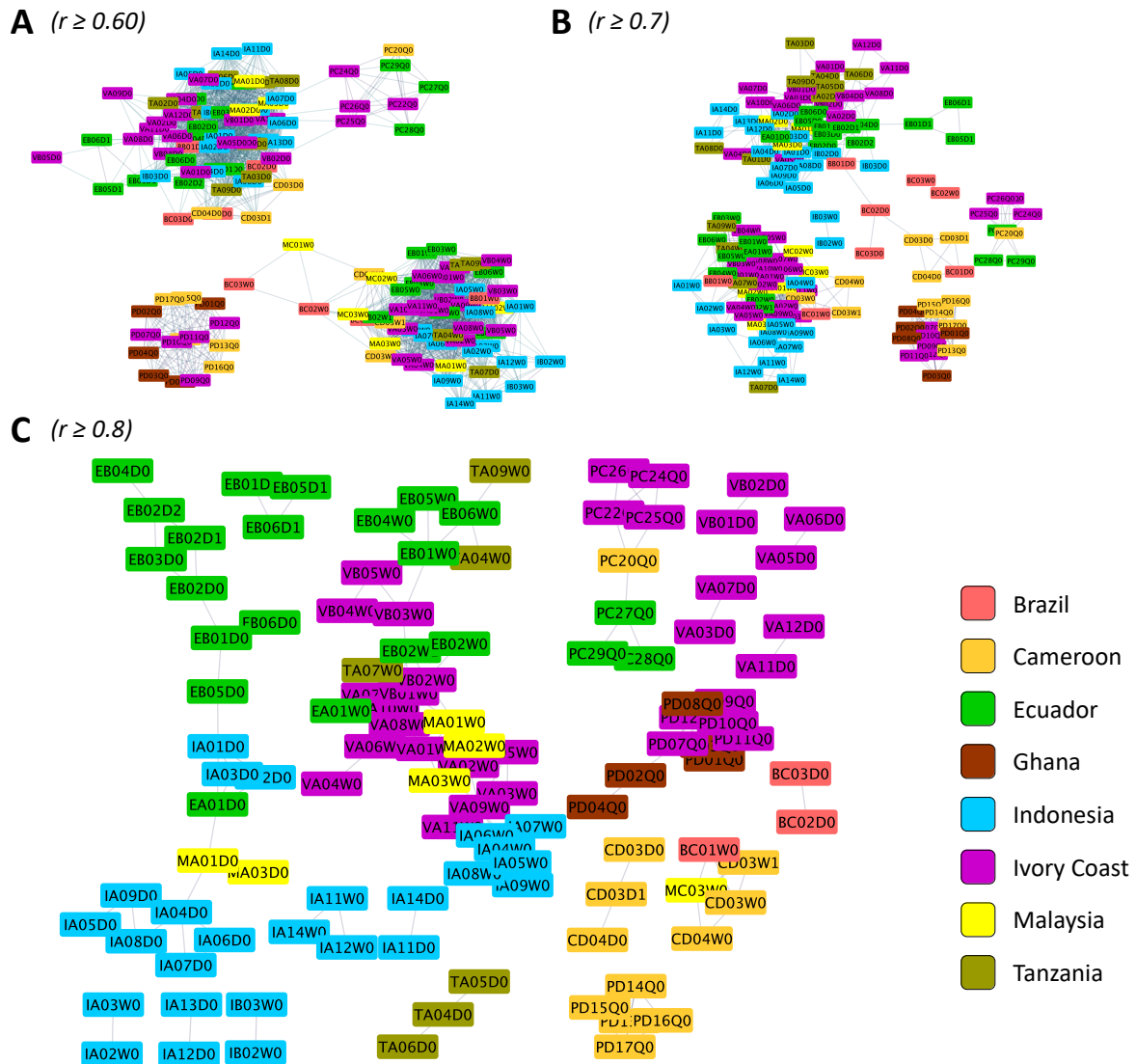rror! Reference source not found.**. See supplementary information for a movie on evolving network as the correlation threshold is progressively increased.

247    **3.4 Country enriched modules at high correlation thresholds**

13

248    As the correlation threshold is further increased, the network breaks into various smaller

249    connected components. The resulting individual connected components primarily have the

250    processing-stage sample type. However, there are more than one component that belong to

251    same color or sample type. This reveals the internal structure of the clusters of samples that

252    initially grouped on the basis of their sample types. This additional sub-structure of the network

253    reveals grouping which now is primarily governed by the samples belonging to same country

254    of origin. This is shown in the networks in Figure 3 for correlation thresholds of 0.6, 0.7 and

255    0.8. Panels A and B provide a bird's-eye view at respective thresholds, while panel C gives a

256    detailed view. In contrast to the legend used in previous figures, we now color the nodes on the

257    basis of countries for a quick comprehension of grouping on the basis of countries. The figure

258    with the previous legend scheme is given as Supplementary Information. It can be seen from

259    the figure that same color nodes tend to be present closer together. This feature is visible more

260    in modules of smaller size, but it is also discernible in larger sized modules. We see that as the

261    correlation threshold is further increased, most of the larger size modules break into smaller

262    module, where nodes belonging to the same origin country are increasingly often connected. It

263    should be noted that processing-stage and country of origin are only the major governing

264    factors, on which grouping of samples is based. Other factors such as variety of cocoa hybrid,

265    harvest season, geographical location and landscape of farm in the country etc, can begin to

266    play an important role with increasing correlation threshold. Hence the clustering is not perfect.

267    The other governing factors can potentially lead to finer sub-modular structures in the network.

268    This situation is more likely to be evident at still higher correlation thresholds.

14

**Figure 3 Country modules.** The structure of correlation network of coc[...] their LCMS profile at correlation thresholds of 0.80, 0.85 and 0.90. At these correlation thresholds, several modules with nodes belonging to the same country of origin are revealed. For a quick and better comprehension and unlike the legend of earlier correlation networks, in this figure different countries are represented through a different color. The networks with same thresholds but with previous annotation (i.e. of Figure 1 and Figure 2) is given in Supplementary Information for comparison. See supplementary information for a movie on evolving network as the correlation threshold is progressively increased.
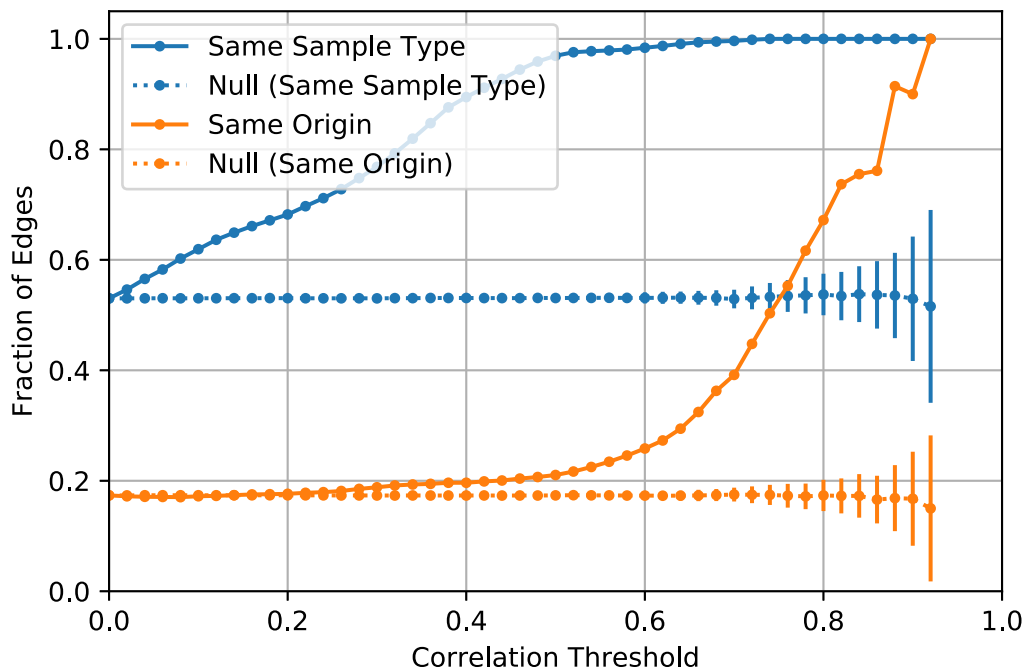
15

280     As the correlation threshold is gradually increased, edges with correlation less than the

281     threshold value are lost from the network. On one hand this leads to increased consideration of

282     the edges with higher correlation in the determination of the layout of the network, while on

283     the other this, naturally, leads to decrease in the number of edges, and when possible, also

284     decreases the number of nodes in the resulting network, resulting in network breakage. The

285     variation of number of edges and number of nodes connecting them is shown in the

286     Supplementary Information file. In our networks here, only the edges greater than the set

287     correlation threshold and corresponding nodes are present. A movie of the network as a

288     function of progressively increasing the threshold is attached as supplementary information

289     which clearly shows the evolving network and separation of samples belonging to different

290     countries.

291     **3.5 Similarity of nodes connected by an edge**

292     As a node in our correlation networks has two attributes, namely the processing-stage sample-

293     type and origin, we define two kinds of similarity for a pair of nodes connected by an edge:

294     sample-type similarity and origin similarity. We define sample-type similarity as the fraction

295     of edges in a network connecting nodes having the same processing-stage sample-type

296     attribute, and origin similarity as the fraction of edges in a network connecting nodes which

297     have same origin attribute. The sample-type and origin similarities as a function of correlation

298     thresholds based on Spearman correlation networks are shown in Figure 4 (solid lines). They

299     differ significantly from each other in terms of both the correlation threshold around which

300     they start to rise and the manner in which they rise. The sample-type similarity starts to increase

301     right from the smaller values of correlation thresholds itself and in a linear manner until it starts

302     to saturate around a correlation threshold value of 0.5 to a similarity value close to 1. This is in

303     agreement with the observed enhancement of the processing-stage sample type character of the

16

304    network architecture right from the beginning of starting values of correlation threshold, to the

305    almost full appearance of processing-stage sample type character at intermediate correlation

306    threshold in large and small connected components (cf. Figure 2). The origin similarity remains

307    almost constant and close to that of null model networks (orange dashed line) for a long range

308    of correlation threshold (up to 0.5) suggesting a weak or almost negligible role in the clustering

309    of nodes belonging to the same origin in the layout of network. Only when the correlation

310    threshold is around 0.5, origin similarity starts to increase, suggesting this is the value of

311    correlation threshold at which the contribution of origin effects start to contribute in clustering

312    of nodes belonging to same origin begins.  This clearly shows that the processing-stage sample

313    type effect precedes the country effects, and the country effects are finer than the sample-type

314    effect. The origin similarity increases exponentially and reaches a value close to 1. This implies

315    that at higher threshold almost all edges connect nodes having same sample type and same

316    country of origin.



317

318   **Figure 4 Connected nodes' similarity.** The sample-type similarity (blue line) starts to increase

319       linearly right from smaller correlation threshold values, reaches close 1 around a correlation

320       threshold value of 0.5. The origin similarity remains constant for a long range of correlation

321       threshold (0, 0.50) and then increases exponentially. The dashed lines show corresponding

322       similarities as expected from an ensemble of control networks.

323   The dashed lines along with error bars show similarity values and standard deviation expected

324   from an ensemble of null model networks (control networks) obtained by randomizing edge

325   weights in the original network (see 2.4 Null model network or control network). The

326   difference between the similarity values from original network and that obtained null model

327   networks points to the fact that the networks at higher correlation thresholds are enriched in

328   edges that have high sample-type and origin similarity. The result corresponding to correlation

329   network generated using Pearson correlation coefficient is given in Supplementary Information

330   file. Both show similar behavior, although at slightly different correlation threshold value.
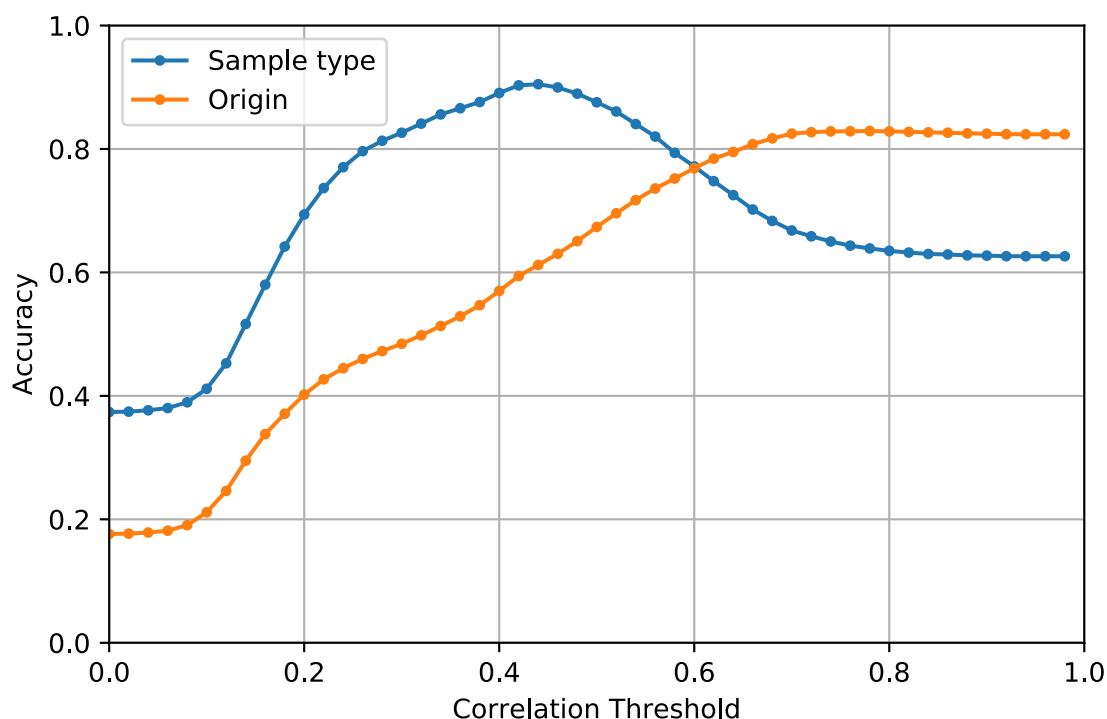
331   **3.6 Closeness of thresholded networks to ideal networks**

332   In this section, we quantify as a function of correlation threshold how accurately our networks

333   represent the expected ideal networks of cocoa samples given their processing-stage sample

334   types or country of origin. We consider two ideal networks, one each for the processing-stage

335   sample type and country of origin. An ideal processing-stage sample type based network will

336   have a link between a pair of its nodes only when both the nodes belong to the same processing-

337   stage sample type, otherwise the link would be absent. Similarly, an ideal origin-based network

338   will have a link between a pair of its nodes only when both the nodes belong to the same

339   country of origin. Thus, in an ideal network based upon processing-stage sample type or

340   country of origin a link is present only between nodes belonging to same sample type, or nodes

341   belonging to same origin, otherwise there is no link between dissimilar nodes. After defining

342   these ideal or true networks, we identify 'true positive' and 'true negative' links by comparing

18

343    the links in the original network at a given correlation threshold (or thresholded network, for

344    short) with the links in the ideal networks. A link is counted to be 'true positive' when the link

345    is present both in the original network at the given threshold and the corresponding ideal

346    network. A link is counted as 'true negative' when the link is absent both in the network at the

347    given threshold and the corresponding ideal network. On the other hand, a link is defined as

348    'false positive' when it is present in the thresholded network but not in the corresponding true

349    network, and 'false negative' when it is absent in the thresholded network but present the true

350    network. An illustration of this scheme through a toy network is provided in Supplementary

351    Information file. Using these terms, we define accuracy $\alpha$ as the fraction of 'true positive' and

352    'true negative' links in an original thresholded network. Accuracy quantifies how close a

353    thresholded network is to the ideal expected network.

354    We find that with increasing correlation thresholds the network becomes closer to the expected

355    true network as demonstrated by increasing values of accuracy for both processing-stage

356    sample type and country of origin Figure 5 (Spearman correlation network; Pearson correlation

357    case in Supplementary Information file). Further, in the region of low correlation threshold the

358    character of the network is closer to that of the expected true network for the processing-stage

359    sample type attribute, and in the region of higher correlation threshold the character of the

360    network is closer to that of the expected true network for country of origin attribute. This result

361    is in agreement with the previous results with formation of processing-stage sample type

362    clusters at lower and intermediate correlation thresholds and of country-based clusters at high

363    correlation thresholds.

364

19

365

**Figure 5 Accuracy of links in thresholded correlation networks, or closeness of a threholded correlation network to expected ideal network based on sample type or origin attributes of cocoa samples.** As the correlation threshold increases the threshold networks become closer to their ideal counterparts. In regions of lower correlation threshold, the thresholded networks are describe more the sample type character of the network than the origin type character. In regions of higher correlation threshold, opposite is true and the thresholded networks are closer in their character to the origin attribute of LC-MS samples. This is coherent with the network pictures at various threshold seen in earlier figures.

## 4. Conclusions and Discussion

We have introduced a new approach for studying grouping in cocoa samples using their LC-MS profile. This new approach is often called 'network science', and it already benefits a multitude of scientific disciplines. Few cases also exist where network approach has been successfully applied in food science for different purposes (Hochberg et al., 2013; Ursem et

20

379 al., 2008; Wang et al., 2017), however, to the best of our knowledge, we apply it for the first

380 time to study the classification of cocoa samples based upon their LC-MS profiles.

381 Classification of cocoa samples on the basis of their country of origin has been found

382 challenging with limited success obtained in cases with the number of countries being few or

383 the origin being on continental scale. Differences in unfermented and fermented samples can

384 be easily seen by simply finding the Spearman correlation between the cocoa samples using

385 their LC-MS profiles (cf. Figure 1). The liquor samples are closer to the fermented samples.

386 However, differentiation on the level of country of origin is only revealed upon further analysis.

387 We make a correlation network using the correlation matrix for cocoa samples, and show that

388 systematic variation of a single parameter, namely correlation threshold, can be used to reveal

389 grouping of cocoa samples on the basis of processing-stage, viz. unfermented, fermented and

390 liquor, and country of origin. In the low and intermediate ranges of correlation threshold

391 processing-stage sample type clusters are revealed, and in the higher range of correlation

392 threshold the clustering of cocoa samples on the basis of country of origin is witnessed. We

393 present our results both qualitatively (cf. Figure 2 and Figure 3) and quantitatively (cf. Figure

394 4 and Figure 5). Besides a successful working approach, our work shows that differentiation

395 of cocoa samples on the level of country of origin is on a more subtle level than their

396 differentiation on the basis of processing-stage sample types.

397 It is worth comparing our approach to an often-used method in similar situation—the principal

398 component analysis (PCA). PCA projects the samples into a lower dimensional space whose

399 axes represent highest possible variation on the basis of the features in the dataset used in the

400 analysis. Often it turns out that this analysis is also able to provide us a view in which samples

401 with different classes well separated. However, there is no binding reason for it to be so, as

402 PCA focuses on maximizing variation amongst the samples on the basis of their features and

21

403   not clustering them per se. Further, only truncated amount of information can be used to

404   visualize the samples as we are limited to a maximum of three dimensions. On the other hand,

405   in a correlation network information from all features (compounds used to calculate

406   correlation) is present. Further, one is able to look at the structure of the network at the level

407   of different amount of information by pruning the network thereby keeping low/high

408   correlations as per need. In this sense, the approach of correlation networks is more

409   sophisticated than that of PCA, omitting PCAs basic philosophy of data reduction.

410   Our study takes into consideration two factors on which cocoa samples may primarily differ:

411   processing-stage and country of origin. However, it is worth noting these are not the only

412   governing factors that affects similarity of cocoa samples. Many other factors such as variety

413   of cocoa hybrid, soil, climate, terrain, harvesting season, farming practices etc. also have

414   significant effects (Acierno et al., 2016; Adeniyi et al., 2019; Arévalo-Hernández et al., 2019;

415   Ehiakpor et al., 2016; Kongor et al., 2016). It would be interesting to consider some of these

416   factors in future works and see in what range of correlation threshold these effects start to

417   matter, or can the inclusion of these additional factors give more clear modules of cocoa

418   samples. Besides providing a new approach to study similarity in cocoa samples, our approach

419   can be a compliment to the traditional approaches in this field.

420

22

428  **References**

429  Acierno, V., Yener, S., Alewijn, M., Biasioli, F., and van Ruth, S. (2016). Factors contributing

430  to the variation in the volatile composition of chocolate: Botanical and geographical origins of

431  the cocoa beans, and brand-related formulation and processing. Food Research International

432  *84*, 86–95.

433  Acierno, V., Alewijn, M., Zomer, P., and van Ruth, S.M. (2018). Making cocoa origin

434  traceable: Fingerprints of chocolates using Flow Infusion - Electro Spray Ionization - Mass

435  Spectrometry. Food Control *85*, 245–252.

436  Adeniyi, S.A., de Clercq, W.P., and van Niekerk, A. (2019). Assessing the relationship between

437  soil quality parameters of Nigerian alfisols and cocoa yield. Agroforest Syst *93*, 1235–1250.

438  Ahn, Y.-Y., Ahnert, S.E., Bagrow, J.P., and Barabási, A.-L. (2011). Flavor network and the

439  principles of food pairing. Sci Rep *1*, 1–7.

440  Albert, R., and Barabási, A.-L. (2002). Statistical mechanics of complex networks. Rev. Mod.

441  Phys. *74*, 47–97.

442  Arévalo-Hernández, C.O., da Conceição Pinto, F., de Souza Júnior, J.O., de Queiroz Paiva, A.,

443  and Baligar, V.C. (2019). Variability and correlation of physical attributes of soils cultivated

444  with cacao trees in two climate zones in Southern Bahia, Brazil. Agroforest Syst *93*, 793–802.

445  Barabási, A.-L., Gulbahce, N., and Loscalzo, J. (2011). Network medicine: a network-based

446  approach to human disease. Nature Reviews Genetics *12*, 56–68.

23

447    Batushansky, A., Toubiana, D., and Fait, A. (2016). Correlation-Based Network Generation,

448    Visualization, and Analysis as a Powerful Tool in Biological Studies: A Case Study in Cancer

449    Cell Metabolism. BioMed Research International.

450    Becker, T., Beber, M.E., Windt, K., and Hütt, M.-T. (2012). The impact of network

451    connectivity on performance in production logistic networks. CIRP Journal of Manufacturing

452    Science and Technology 5, 309–318.

453    Bertoldi, D., Barbero, A., Camin, F., Caligiani, A., and Larcher, R. (2016). Multielemental

454    fingerprinting and geographic traceability of Theobroma cacao beans and cocoa products. Food

455    Control 65, 46–53.

456    Borgatti, S.P., Mehra, A., Brass, D.J., and Labianca, G. (2009). Network Analysis in the Social

457    Sciences. Science 323, 892–895.

458    Caligiani, A., Palla, L., Acquotti, D., Marseglia, A., and Palla, G. (2014). Application of 1H

459    NMR for the characterisation of cocoa beans of different geographical origins and fermentation

460    levels. Food Chemistry 157, 94–99.

461    Claussen, J.C., Skiecevičienė, J., Wang, J., Rausch, P., Karlsen, T.H., Lieb, W., Baines, J.F.,

462    Franke, A., and Hütt, M.-T. (2017). Boolean analysis reveals systematic interactions among

463    low-abundance species in the human gut microbiome. PLOS Computational Biology 13,

464    e1005361.

465    Dogrusoz, U., Giral, E., Cetintas, A., Civril, A., and Demir, E. (2009). A Layout Algorithm for

466    Undirected Compound Graphs. Inf. Sci. 179, 980–994.

24

467    D'Souza, R.N., Grimbs, S., Behrends, B., Bernaert, H., Ullrich, M.S., and Kuhnert, N. (2017).

468    Origin-based polyphenolic fingerprinting of Theobroma cacao in unfermented and fermented

469    beans. Food Research International *99*, 550–559.

470    Ehiakpor, D.S., Danso-Abbeam, G., and Baah, J.E. (2016). Cocoa farmer's perception on

471    climate variability and its effects on adaptation strategies in the Suaman district of western

472    region, Ghana. Cogent Food & Agriculture *2*, 1210557.

473    Fayeulle, N., Meudec, E., Boulet, J.C., Vallverdu-Queralt, A., Hue, C., Boulanger, R.,

474    Cheynier, V., and Sommerer, N. (2019). Fast Discrimination of Chocolate Quality Based on

475    Average-Mass-Spectra Fingerprints of Cocoa Polyphenols. J. Agric. Food Chem. *67*, 2723–

476    2731.

477    Grimbs, A., Klosik, D.F., Bornholdt, S., and Hütt, M.-T. (2019). A system-wide network

478    reconstruction of gene regulation and metabolism in Escherichia coli. PLOS Computational

479    Biology *15*, e1006962.

480    Guehi, T.S., Zahouli, I.B., Ban-Koffi, L., Fae, M.A., and Nemlin, J.G. (2010). Performance of

481    different drying methods and their effects on the chemical quality attributes of raw cocoa

482    material. International Journal of Food Science & Technology *45*, 1564–1571.

483    Hagberg, A.A., Schult, D.A., and Swart, P.J. (2008). Exploring Network Structure, Dynamics,

484    and Function using NetworkX. In Proceedings of the 7th Python in Science Conference, p.

485    Hochberg, U., Degu, A., Toubiana, D., Gendler, T., Nikoloski, Z., Rachmilevitch, S., and Fait,

486    A. (2013). Metabolite profiling and network analysis reveal coordinated changes in grapevine

487    water stress response. BMC Plant Biology *13*, 184.

25

488    Jeong, H., Tombor, B., Albert, R., Oltvai, Z.N., and Barabási, A.-L. (2000). The large-scale

489    organization of metabolic networks. Nature *407*, 651–654.

490    Kamada, T., and Kawai, S. (1989). An algorithm for drawing general undirected graphs.

491    Information Processing Letters *31*, 7–15.

492    Kongor, J.E., Hinneh, M., de Walle, D.V., Afoakwa, E.O., Boeckx, P., and Dewettinck, K.

493    (2016). Factors influencing quality variation in cocoa (Theobroma cacao) bean flavour profile

494    — A review. Food Research International *82*, 44–52.

495    Kuhnert, N., Dairpoosh, F., Yassin, G., Golon, A., and Jaiswal, R. (2013). What is under the

496    hump? Mass spectrometry based analysis of complex mixtures in processed food – lessons

497    from the characterisation of black tea thearubigins, coffee melanoidines and caramel. Food

498    Funct. *4*, 1130–1147.

499    Kumar, S., and Deo, N. (2012). Correlation and network analysis of global financial indices.

500    Phys. Rev. E *86*, 026101.

501    Kumar, S., Mahajan, S., and Jain, S. (2018). Feedbacks from the metabolic network to the

502    genetic network reveal regulatory modules in E. coli and B. subtilis. PLOS ONE *13*, e0203311.

503    Kumari, N., Grimbs, A., D'Souza, R.N., Verma, S.K., Corno, M., Kuhnert, N., and Ullrich,

504    M.S. (2018). Origin and varietal based proteomic and peptidomic fingerprinting of Theobroma

505    cacao in non-fermented and fermented cocoa beans. Food Research International *111*, 137–

506    147.

507    Lima, L.J.R., Almeida, M.H., Nout, M.J.R., and Zwietering, M.H. (2011). Theobroma cacao

508    L., "The Food of the Gods": Quality Determinants of Commercial Cocoa Beans, with Particular

26

509    Reference to the Impact of Fermentation. Critical Reviews in Food Science and Nutrition *51*,

510    731–761.

511    Marseglia, A., Acquotti, D., Consonni, R., Cagliani, L.R., Palla, G., and Caligiani, A. (2016).

512    HR MAS 1H NMR and chemometrics as useful tool to assess the geographical origin of cocoa

513    beans – Comparison with HR 1H NMR. Food Research International *85*, 273–281.

514    McKinney, W. (2010). Data Structures for Statistical Computing in Python. In Proceedings of

515    the 9th Python in Science Conference, S. van der Walt, and J. Millman, eds. pp. 51–56.

516    McKinney, W. (2011). pandas: a Foundational Python Library for Data Analysis and Statistics

517    | R (Programming Language) | Database Index. In PyHPC 2011, p.

518    Megías-Pérez, R., Grimbs, S., D'Souza, R.N., Bernaert, H., and Kuhnert, N. (2018). Profiling,

519    quantification and classification of cocoa beans based on chemometric analysis of

520    carbohydrates using hydrophilic interaction liquid chromatography coupled to mass

521    spectrometry. Food Chemistry *258*, 284–294.

522    Milev, B.P., Patras, M.A., Dittmar, T., Vrancken, G., and Kuhnert, N. (2014). Fourier

523    transform ion cyclotron resonance mass spectrometrical analysis of raw fermented cocoa beans

524    of Cameroon and Ivory Coast origin. Food Research International *64*, 958–961.

525    Namaki, A., Shirazi, A.H., Raei, R., and Jafari, G.R. (2011). Network analysis of a financial

526    market based on genuine correlation and threshold method. Physica A: Statistical Mechanics

527    and Its Applications *390*, 3835–3841.

528    Newman, M. (2003). The Structure and Function of Complex Networks. SIAM Rev. *45*, 167–

529    256.

27

530 Oberrauter, L.-M., Januszewska, R., Schlich, P., and Majchrzak, D. (2018). Sensory evaluation
531 of dark origin and non-origin chocolates applying Temporal Dominance of Sensations (TDS).
532 Food Research International *111*, 39–49.

533 Oliveira, L.F., Braga, S.C.G.N., Augusto, F., Hashimoto, J.C., Efraim, P., and Poppi, R.J.
534 (2016). Differentiation of cocoa nibs from distinct origins using comprehensive two-
535 dimensional gas chromatography and multivariate analysis. Food Research International *90*,
536 133–138.

537 Ozretic-Dosen, D., Skare, V., and Krupka, Z. (2007). Assessments of country of origin and
538 brand cues in evaluating a Croatian, western and eastern European food product. Journal of
539 Business Research *60*, 130–136.

540 Pluskal, T., Castillo, S., Villar-Briones, A., and Orešič, M. (2010). MZmine 2: Modular
541 framework for processing, visualizing, and analyzing mass spectrometry-based molecular
542 profile data. BMC Bioinformatics *11*, 395.

543 Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N.,
544 Schwikowski, B., and Ideker, T. (2003). Cytoscape: A Software Environment for Integrated
545 Models of Biomolecular Interaction Networks. Genome Res. *13*, 2498–2504.

546 Sirbu, D., Grimbs, A., Corno, M., Ullrich, M.S., and Kuhnert, N. (2018). Variation of
547 triacylglycerol profiles in unfermented and dried fermented cocoa beans of different origins.
548 Food Research International *111*, 361–370.

549 Szklarczyk, D., Franceschini, A., Wyder, S., Forslund, K., Heller, D., Huerta-Cepas, J.,
550 Simonovic, M., Roth, A., Santos, A., Tsafou, K.P., et al. (2015). STRING v10: protein–protein
551 interaction networks, integrated over the tree of life. Nucleic Acids Res *43*, D447–D452.

28

552    Ursem, R., Tikunov, Y., Bovy, A., van Berloo, R., and van Eeuwijk, F. (2008). A correlation

553    network approach to metabolic data analysis for tomato fruits. Euphytica *161*, 181.

554    Vázquez-Ovando, A., Molina-Freaner, F., Nuñez-Farfán, J., Betancur-Ancona, D., and

555    Salvador-Figueroa, M. (2015). Classification of cacao beans (Theobroma cacao L.) of southern

556    Mexico based on chemometric analysis with multivariate approach. Eur Food Res Technol

557    *240*, 1117–1128.

558    Wang, L., Sun, X., Weiszmann, J., and Weckwerth, W. (2017). System-Level and Granger

559    Network Analysis of Integrated Proteomic and Metabolomic Dynamics Identifies Key Points

560    of Grape Berry Development at the Interface of Primary and Secondary Metabolism. Front.

561    Plant Sci. *8*.

562    Wold, S., Esbensen, K., and Geladi, P. (1987). Principal component analysis. Chemometrics

563    and Intelligent Laboratory Systems *2*, 37–52.

564

565

566

29