# Near-infrared hyperspectral imaging for non-destructive classification of commercial tea products

Puneet Mishra[1], Alison Nordon[1], Julius Tschannerl[2], Guoping Lian[3,4], Sally Redfern[3], Stephen Marshall[2]

[1]*WestCHEM, Department of Pure and Applied Chemistry and Centre for Process Analytics and Control Technology, University of Strathclyde, 295 Cathedral Street, Glasgow, G1 1XL, United Kingdom*

[2]*Hyperspectral Imaging Centre, Department of Electronic and Electrical Engineering, University of Strathclyde, 204 George Street, Glasgow, G1 1XW, United Kingdom*

[3]*Unilever R&D Colworth, Colworth House, Sharnbrook, Bedford MK44 1LQ, United Kingdom*

[4]*Department of Chemical and Process Engineering, University of Surrey, Guildford GU2 7XH, United Kingdom*

*Corresponding authors email: puneet.mishra@strath.ac.uk , alison.nordon@strath.ac.uk*

**Abstract**

Tea is the most consumed manufactured drink in the world. In recent years, various high end analytical techniques such as high-performance liquid chromatography have been used to analyse tea products. However, these techniques require complex sample preparation, are time consuming, expensive and require a skilled analyst to carry out the experiments. Therefore, to support rapid

19    and non-destructive assessment of tea products, the use of near infrared (NIR) (950-1760 nm)

20    hyperspectral imaging (HSI) for classification of six different commercial tea products (oolong,

21    green, yellow, white, black and Pu-erh) is presented. To visualise the HSI data, linear (principal

22    component analysis (PCA) and multidimensional scaling (MDS)) and non-linear (t-distributed

23    stochastic neighbour embedding (t-SNE) and isometric mapping (ISOMAP)) data visualisation

24    methods were compared. t-SNE provided separation of the six commercial tea products into three

25    groups based on the extent of processing: minimally processed, oxidised and fermented. To

26    perform the classification of different tea products, a multi-class error-correcting output code

27    (ECOC) model containing support vector machine (SVM) binary learners was developed. The

28    classification model was further used to predict classes for pixels in the HSI hypercube to obtain

29    the classification maps. The SVM-ECOC model provided a classification accuracy of 97.41±0.16

30    % for the six commercial tea products. The methodology developed provides a means for rapid,

31    non-destructive, *in situ* testing of tea products, which would be of considerable benefit for process

32    monitoring, quality control, authenticity and adulteration detection.

33    ***Keywords:*** Imaging spectroscopy, hypercube, multivariate, data visualisation, neighbourhood

34    methods

35    **1. Introduction**

36    Being the oldest beverage, tea is the most consumed drink in the world (Sang, 2016). Different tea

37    products exist due to different processes for freshly harvested tea leaves (Lv et al., 2013). There

38    are six main types of tea products, i.e. oolong, green, yellow, white, black and Pu-erh (Chang,

39    2015), which differ in terms of processing (see Figure 1). Green, yellow and white tea products

40    undergo minimal processing, oolong and black tea products have been oxidised while Pu-erh tea

41    has been fermented. The chemical composition of fresh tea (*Camellia sinensis*) leaves is a complex

42    mixture of caffeine, polyphenols, polysaccharides and nutrients such as protein, amino acids,

43    lipids, and vitamins (Ruan et al., 2010). Typically, during the processing of fresh tea leaves, such

44    as oxidation and fermentation, they undergo chemical compositional changes. Free amino acids,

45    total tea polyphenols, soluble sugars, and caffeine are the four major chemical components that

46    determine the nature and quality of the final tea products (Ozturk et al., 2016).

47    *Figure 1: Processing steps for different tea products starting from fresh green tea leaves to final*
48    *products.*

49
50    Analytical methods used to measure chemical constituents as quality indicators of plant-based

51    products include high-performance liquid chromatography (HPLC) (Nieh et al., 2009), liquid

52    chromatography/mass spectrometry (LC/MS) (Tan et al., 2016), gas chromatography/mass

53    spectrometry (GC/MS) (Jing et al., 2017) and electrochemical systems (Kumar et al., 2016)

54    (Domínguez et al., 2015). However, these methods have complex sample preparation, are time

55    consuming, expensive and require a skilled analyst to carry out the experiments (Li et al., 2017).

56    A non-destructive technique that has been used for analysis of tea processes and quality monitoring

57    is e-nose (Yaroshenko et al., 2014) (Sharma et al., 2015). E-nose devices usually include an array

58    of metal oxide sensors which respond to the amount of biochemical volatiles coming into contact

59    with the corresponding sensor surface to explain the chemical profile (Bhattacharyya et al., 2007).

60    However, a major disadvantage of e-nose sensors is that they are affected by environmental

61    conditions such as temperature and humidity, which leads to sensor drift (Baldwin et al., 2011).

62    In recent years, there has been increasing interest in the use of optical spectroscopic techniques for

63    rapid, non-destructive assessment of food products. NIR spectroscopy is particularly attractive for

64    this purpose, where changes in the NIR spectral profiles can be correlated to perform qualitative

65   and quantitative analysis of food products (Qu et al., 2015; Fu and Ying, 2016). NIR spectroscopy

66   has been explored for discrimination (He et al., 2007; Chen et al., 2009), identification (Chen et

67   al., 2007; Wang et al., 2015) and quality assessment (Panigrahi et al., 2016) of tea products. Also

68   reported for non-destructive tea analysis are emerging studies utilising imaging techniques for the

69   identification (Chen et al., 2008), classification (Wang et al., 2015) and for evaluation of sensory

70   quality (Zhu et al., 2017) of tea products. Integration of spectroscopy and imaging is known as

71   hyperspectral imaging (HSI) and use of NIR-HSI still seems unexplored in its application to the

72   analysis of tea products.

73   HSI has been widely used in remote sensing for military applications (Goetz et al., 1985), but it is

74   now popular in scientific domains such as forensics (Edelman et al., 2012), medical (Lu et al.,

75   2014), food (Pu et al., 2015), pharmaceutical (Kandpal et al., 2016) and plants (Mishra et al., 2017).

76   There are reports of the use of HSI for the understanding of different food products such as coffee

77   (Nansen et al., 2016), tobacco (Garcia-Allende et al., 2008), and seeds of vegetable and fruits

78   (Shrestha et al., 2016; Kandpal et al., 2016). Some applications of HSI of tea have been reported

79   but these studies only considered a single variety of tea and measured the visible and very near

80   infrared (VNIR) range (around 400-1000 nm), which is dominated by the pigments and physical

81   characteristics of the samples (Zhao et al., 2009; Xie et al., 2015). In comparison to the VNIR

82   region, the NIR region provides more detailed chemical information such as overtones resulting

83   from the molecular vibration of O-H, C-H, N-H bonds and their combinations, which can support

84   a better classification system based on the chemistry of the samples (Mishra et al., 2016).

85   The aim of the present work is to demonstrate the use of NIR (950-1760 nm) HSI for rapid, non-

86   destructive classification of six different commercial tea products (oolong, green, yellow, white,

87   black and Pu-erh). The study investigates and compares four different dimensionality reduction

88     techniques (linear and non-linear) to visualise the high dimensional HSI tea data. Furthermore,

89     multi-class support vector machine (SVM) modelling has been performed to generate spatial

90     classification maps of tea products.

## 2. Materials and Methods

91

### 2.1. Samples

92

93     Six commercial tea samples were obtained from the local market (Glasgow, United Kingdom).

94     The samples were obtained in airtight sealed packaging and stored at ambient temperature. All

95     samples of tea were in loose-leaf form. Black, green and white tea were from Vahdam Teas (New

96     Delhi, India), oolong tea was from Yamamotoyama (California, USA), Pu-erh tea was from The

97     Tea Makers of London (London, United Kingdom) and yellow tea was of an unspecified Chinese

98     origin. The six tea products can also be broadly grouped as minimally processed (green, white and

99     yellow), oxidised (black and oolong tea) and fermented (Pu-erh tea). The samples for each imaging

100    experiment were transferred on the day of analysis into a black plastic circular container (diameter

101    = 3.3 cm, depth = 1.3 cm). A different cap was used for each tea to avoid any cross-contamination.

### 2.2. Hyperspectral imaging measurements

102

103    Imaging was performed with a push-broom line scan HSI camera (*Model name*: RedEye 1.7) from

104    INNO-SPEC (Nurnberg, Germany). The camera has an InGaAs sensor and generates a spatial map

105    of 320 x 256 pixels in the spectral range of 950 - 1760 nm. The pixel size was 30 x 30 $\mu m^2$ and

106    the spectral resolution was 3.2 nm. The camera communicated with the computer via a gigabit

107    Ethernet connection. The lighting was provided by two halogen light sources 50 W each and the

108    integration time used was 300 ms. Imaging was performed by placing the samples over the

109    translation stage which was controlled by an independent stage motor connected to the computer

110    system (Zolix TSA 200 BF). The speed of the translation stage was optimised before image

111    acquisition to avoid any distortion in the shape of the image arising from the overlapping of the

112    spectral information in the adjacent pixels. The image acquisition and management of settings

113    (integration time) were performed using the software interface called SiCAP provided with the

114    camera by INNO-SPEC. Images were first acquired of six different tea samples placed adjacent to

115    each other in their respective sample containers in the field of view of the camera. An image was

116    then acquired of black, Pu-erh and oolong teas where each tea occupied approximately a third of

117    the volume of the sample container; the teas were not physically mixed. Finally, equal proportions

118    of all six tea samples were mixed, by manually shaking the different tea products in a container,

119    and an image of the mixture was acquired. One image was acquired of each sample, with each

120    image comprising more than 2000 pixels (spectra) for the individual tea samples and more than

121    11200 pixels for the samples containing more than one type of tea. An illustration of the HSI setup

122    configured for imaging of tea samples can be found in Figure 2.

123    *Figure 2: Illustrative diagram for the hyperspectral imaging setup used to acquire the images of tea*

124    *samples.*

125    **2.3. Data analysis**

126    **2.3.1. Pre-processing of HSI data**

127    The data cubes not only contain information about the samples imaged but also consist of different

128    unwanted influences in signal resulting from factors such as illumination intensity, the detector

129    sensitivity and transmission properties of the optics. The effects resulting from these factors are

130    both wavelength dependent and independent. To correct for these effects, radiometric calibration

131    was performed using dark and white reference images acquired along with the samples. The

132    correction was performed for every pixel in the HS image according to equation 1:

$$I_{R(i,j)} = \frac{I_{raw(i,j)} - I_{dark(i,j)}}{I_{white(i,j)} - I_{dark(i,j)}} \qquad (1)$$

133

134    where, $I_R$ is the calibrated reflectance image, $I_{raw}$ is the raw intensity image measured from the test

135    sample, $I_{dark}$ is the intensity of the dark response, $I_{white}$ is the intensity for the uniform white

136    reference and $i$ and $j$ were spatial coordinates over the image.

137    Often, the radiometric correction is sufficient to remove the effects of illumination inhomogeneity

138    from the spectral data, however, when the sample surfaces are not uniform, as in the case of

139    samples of loose tea leaves, the light scattering during diffuse reflection causes additive and

140    multiplicative effects (Mishra et al., 2016). These scattering effects lead to baseline shifts in the

141    spectrum and variation in the global intensity, which is again dependent on the wavelength.

142    Standard normal variate (SNV) is a very common technique used in NIR spectroscopy to remove

143    these effects (Barnes et al., 1989). In SNV, the mean and standard deviation of each spectrum for

144    each pixel are calculated, the mean is subtracted, and the standard deviation is used to normalise

145    the difference. This transformation normalises each spectrum to zero mean and unit standard

146    deviation. Before applying the SNV transform, the spectral range was reduced from 950 - 1760

147    nm to 967 nm - 1700 nm, to remove the noisy regions at the edges of the spectral range, and

148    converted to absorbance. Further, the spectral absorbance profiles were smoothed with a Savitzky-

149    Golay filter (15-point width and second order polynomial) (Savitzky and Golay, 1964). The *savgol*

150    and *snv* functions from PLS toolbox (version 8.11, Eigenvector Research Inc., USA) were used.

151    All visualisation and classification analysis was performed on the pre-processed spectra. The pre-

152 processed pure spectra of six pure tea samples were extracted using Matlab's (R2016b,

153 Mathworks, USA) *roipoly* function. The *roipoly* function provides a graphical user interface in

154 Matlab to extract the information from each image over the manually selected locations.

### 155 *2.3.2. Principal Component Analysis*

156 Principal component analysis (PCA) introduced by Pearson in 1901 belongs to the family of linear

157 methods for visualising high dimension data (Wold et al., 1987). In PCA, a set of observations

158 containing correlated variables is orthogonally transformed to linearly uncorrelated variables

159 defined as principal components (PCs). In PCA, the transformation is performed to retain the major

160 amount of variability in the dataset.

161 The PCA decomposition model for a given observation data matrix X can be understood as

162 equation 2:

$$X = TW^T \hspace{4cm} (2)$$

163

164 where T is the score in the lower dimension explained by the number of PCs specified and W is a

165 *p × p (p denotes number of variables)* matrix whose columns are the eigenvectors of $X^T X$.

166 In the case of dimensionality reduction, the aim is to preserve the maximum amount of meaningful

167 variation present in the dataset. The extracted PCs define a new orthonormal basis set which can

168 be used to transform the data from a high dimension space to the lower space explained by the

169 PCs. PCA from a dimensionality reduction perspective can be understood as minimising the

170 squared reconstruction error as given in equation 3.

171 $$min\|TW^T - T_r W_r^T\|^2 \qquad (3)$$

172 where, $TW$ and $T_r W_r$ are the reconstructed original dataset in higher and lower dimensional space

173 respectively. Minimisation of the reconstruction error results in the maximisation of the

174 information that was present in the higher dimensional space when defined in the lower

175 dimensional space given by the significant number of PCs. To interpret the data in two or three

176 dimensional plots, the respective PCs can be selected and used for transformation to the orthogonal

177 axes represented by the PCs. Transformation from a higher dimension to a lower dimension can

178 be performed as in equation 4.

179 $$\hat{X}_r = XW \qquad (4)$$

180 *2.3.3. Multi-Dimensional Scaling*

181 Multi-dimensional scaling (MDS) is a linear method for visualising high dimensional data (Cox et

182 al., 2000). MDS performs a transformation by preserving the between object distances from the

183 higher dimension to lower dimension. The MDS utilises calculation of the Euclidean distances for

184 each data point in the multidimensional space to capture the pattern. The distances are defined as

185 a symmetric distance matrix ($D$). MDS attempts to find data points in a specified (d-dimensional)

186 space such that the Euclidean distance between data points ($\hat{D}$) is similar to the distance in higher

187 dimensional space. The minimisation function can be understood as equation 5:

188 $$min \sum_i \sum_j \|d_{ij} - \hat{d}_{ij}\|^2 \qquad (5)$$

189     where, $D = d_{ij} = \| x_i - x_j \|^2$ and $D\hat{} = d_{ij} = \| y_i - y_j \|^2$ explaining the Euclidean distance between

190     points in high $(x_i, x_j)$ and low dimensional space $(y_i, y_j)$, respectively. $i, j$ denotes specific position

191     of point.

192     *2.3.4. Isometric Mapping*

193     Isometric mapping (ISOMAP) belongs to the family of non-linear techniques for visualising high

194     dimensional data (Tenenbaum, 1998; Balasubramanian and Schwartz, 2002). ISOMAP can be

195     understood as a generalised non-linear form of MDS which utilises the geodesic space accounting

196     for the non-linearity in the high dimensional data manifold. The geodesic distance is defined as

197     the shortest distance between two data points on a curved surface of a non-linear manifold. As a

198     first step, ISOMAP approximates a neighbourhood graph by identifying k nearest neighbours

199     (kNNs) or selecting neighbourhood data points based on any other condition for every data points.

200     The geodesic distance is then approximated for all the pairs of data points on the neighbourhood

201     graph. Finally, the distance data obtained from the graph is embedded to a lower dimension

202     Euclidean space using MDS as shown in equation (6).

$$\min \sum_i \sum_j \|D_G - D_E\|^2 \qquad (6)$$

203

204     where, $D_G$ and $D_E$ explaining the geodesic and Euclidean distance between points in high and low

205     dimensional space, respectively.

206

207     *2.3.5. t-Distributed Stochastic Neighbour Embedding*

208    t-distributed stochastic neighbour embedding (t-SNE) is a non-linear technique used to visualise

209    high dimensional data in two or three dimensional scatter plots (Maaten and HInton, 2008). The

210    main objective of t-SNE is to model the similar points using nearby points (small pairwise

211    distance) and the dissimilar points using distant points (large pairwise distances). As a first step,

212    to represent the similarity, the t-SNE converts high-dimensional Euclidean distances between data

213    points into conditional probabilities using a Gaussian distribution. The joint probability for a data

214    point $x_j$ to $x_i$ can be calculated with equation (7):

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_k \sum_{k \neq l} \exp(-\|x_k - x_l\|^2 / 2\sigma_i^2)}, \qquad (7)$$

215

216    The conditional probability represents the probability that $x_i$ will pick $x_j$ as a neighbour based on

217    the proportion of probability density under a Gaussian centred at $x_i$. If the points are near then the

218    value of $p_{i|j}$ will be higher compare to the points far away. Furthermore, the conditional

219    probabilities are symmetrised to reduce the effects of outliers by setting (8):

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N} \qquad (8)$$

220

221    To represent joint probabilities in the low dimensional map $q_{ij}$, t-SNE utilises a heavy tailed

222    Student t-distribution. The benefit of using a heavy tailed distribution is that it makes the joint

223    probabilities invariant to changes in the scale of the map. The joint probabilities $q_{ij}$ can be

224    estimated by (9):

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k \sum_{k \neq l}(1 + \|y_k - y_l\|^2)^{-1}}, \qquad (9)$$

225

226 Finally, the t-SNE minimises a single Kullback-Leibler (*KL*) divergence between a joint

227 probability distribution, P, in the high-dimensional space and a joint probability distribution, Q, in

228 the low-dimensional space as can be understood from equation (10):

$$KL(P\|Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}} \qquad (10)$$

229

230 The minimisation of the KL divergence is performed using a gradient descent algorithm with

231 respect to the locations of the points in the map $y_i$.

232 All the data visualisation methods (PCA, MDS, ISOMAP and t-SNE) were implemented in Matlab

233 using the Toolbox for Dimensionality Reduction (https://lvdmaaten.github.io/drtoolbox/) (Maaten

234 et al., 2009; Maaten and Hinton, 2008). The Mahalanobis distance (Mahalanobis, 1936) was used

235 to assess the separation of the clusters identified with the different data visualisation methods.

236 ***2.3.6. Support vector machines for multi-class classification***

237 Support vector machines (SVMs) are supervised non-probabilistic learning models which utilise

238 hyperplanes to define the decision boundaries for performing classification (Vapnik and Vapnik,

239 1998). The SVM algorithms are usually developed to perform a binary classification, however,

240 SVM can be used for multi-class classification problems by utilising several independent binary

241 classifiers. This can be performed by combining it with ensemble methods such as error correcting

242 output codes (ECOC). The ECOC deals with the multi-class classification problem by converting

243     it into several independent binary classification problems. A wide range of applications of SVM

244     to process HS images can be seen (García Allende et al., 2008) (Mountrakis et al., 2011).

245     In the present work, the ECOC-SVM algorithm available in Matlab's Statistics and Machine

246     Learning Toolbox (R2016b) was implemented to perform the classification utilising the

247     classification learner application. ECOC-SVM uses a one-versus-all coding design, in which for

248     each binary learner one class is assigned a positive value and all others are assigned negative

249     values. To map the data to the higher dimension, a radial basis function (RBF) kernel (scale

250     parameter=10) was used. The RBF kernel has the benefit of non-linearly mapping the sample to

251     the higher dimensional space for dealing with a non-linear relationship between observation and

252     classes. For every pure tea sample, spectra (967 - 1700 nm) were extracted from 200 pixels, which

253     were selected at random from the image collected, leading to 1200 spectra in total for calibration

254     of the classification model. Validation of the model was performed with a 10-fold cross-validation

255     method. Furthermore, to have confidence in the model accuracy, the model was recalibrated with

256     1200 iterations and the mean and standard deviation were noted. The trained classifier was further

257     used to generate the classification maps of the HS images. The HSI cubes were first unfolded from

258     a 3D map ($n \times p \times k$) to a 2D matrix ($np \times k$) and then the class of every row of the matrix

259     (representing the pixel) was predicted, where n, p, k defined the x, y and z dimension of data. After

260     prediction, the matrix ($np \times 1$) was reshaped to the original image dimension ($n \times p$).

261     **3. Results**

262     ***3.1. Spectral profiles of tea samples***

263

264 *Figure 3: Absorbance spectra of pure tea samples of yellow, oolong, green, black, white and Pu-*
265 *erh. (a). Mean absorbance spectra (n = 200). (b) Mean spectra after pre-processing (SNV and*
266 *Savitzky- Golay smoothing), and (c) standard deviation of the absorbance spectra and spectra*
267 *after pre- processing. The vertical green lines denote the positions of the main peaks.*

268

269   Figure 3 presents the spectral profiles of individual tea samples. Figure 3(a) presents the mean

270   absorbance spectra calculated from the 200 spectra extracted for each of the six tea samples

271   (yellow, oolong, green, black, white and Pu-erh), Figure 3(b) presents the mean spectra after pre-

272   processing with Savitzky-Golay filtering followed by SNV, and Figure 3(c) presents the standard

273   deviation of the spectra before and after pre-processing. From Figure 3(a), it can be seen that the

274   absorbance spectra of different tea samples contain scattering effects leading to baseline shifts.

275   These effects can also be seen in the standard deviation plot in Figure 3(c) for the absorbance

276   spectra (red), where the standard deviation over the entire spectral range is approximately constant.

277   These scattering effects can bias modelling of the data, therefore, they were removed via pre-

278   processing. In Figure 3(b), it can be seen that after pre-processing, differences in spectra at various

279   wavelengths have emerged, and so spectral differences corresponding to different teas can be

280   noted. Scattering effects arise in the imaging experiments as the inhomogeneity in the size of the

281   loose leaves does not get compensated for by the flat surface of the white reflectance standard used

282   for radiometric calibration.

283   In Figure 3(c), it can be noted that the pre-processing reveals the spectral variation arising from

284   differences in the tea, which was previously dominated by the effects of light scattering. In Figure

285   3(b), various peaks (depicted by the green vertical lines) can be identified at representative

286   wavelengths. In previous works, the peaks at 1131, 1654 and 1666 nm were found to be

287    representative of the total tea polyphenols (Chen et al., 2006; Bian et al., 2010; Bian et al., 2013),

288    1361 nm is representative of moisture content (Panigrahi et al., 2016), 1093-1121 nm for

289    thearubigin components of TRS1 (Panigrahi et al., 2016), 1492 nm corresponds to free amino acids

290    (Bian et al., 2010), 1176 nm is a second overtone C-H (Tan et al., 2012) and 1390 nm for the $CH_2$-

291    overtone (Lee et al., 2014).

292    ***3.2. Visualising high dimensional data***

293

294    *Figure 4: 2-Dimensional scatter plots for visualising high dimensional tea data. (a). Principal Component*

295    *Analysis (PCA), (b). Multidimensional Scaling (MDS), (c). Isometric Mapping (ISOMAP), and (d). t-*

296    *distributed Stochastic Neighbour Embedding (t-SNE). In all the plots, the first dimension is represented in*

297    *the x-axis and the second in the y-axis, and the six tea products are coloured as follows: Pu-erh (pink),*

298    *black (sky blue), oolong (yellow), green (green), white (blue) and yellow (red).*

299    To visualise the high dimensional data in the lower dimension, the 256-dimensional HSI data were

300    transformed to 2-dimensional plots using PCA, MDS, ISOMAP and t-SNE as shown in Figure 4.

301    It can be seen clearly in Figure 4 that the t-SNE (Figure 4(d)) outperforms PCA, MDS and

302    ISOMAP (Figures 4(a), 4(b) and 4(c), respectively) regarding identification of the maximum

303    number of separate clusters. These separate clusters correspond to different tea products and their

304    representation as separate clusters in the plots signifies that the visualisation method is able to

305    preserve the structure of the data on transformation from a high dimensional space to a lower

306    dimensional space. In general, all the methods were able to separate the Pu-erh tea (pink) from all

307    other tea samples. The reason for this can be seen in Figure 3(b) where Pu-erh tea (sky blue) has a

308    very different spectral signature compared to the other tea samples. This is likely to be because the

309    Pu-erh tea undergoes very different processing, which includes microbial fermentation of sun dried

310    leaves (Lv et al., 2013), compared to the other teas.

311    It can be seen in Figure 4((a), (b) and (c)) that with the exception of Pu-erh tea, all other types of

312    tea samples are mixed and their clear distinction is not possible. In comparison, black and oolong

313    tea are identified as separate clusters with t-SNE. However, while t-SNE was not able to separate

314    the green, yellow and white tea, it still provided better separation of these three teas as shown in

315    Figure 4(d). Green, yellow and white teas appear in the same cluster as they have similar spectral

316    signatures (see Figure 3(b)). This may arise from the fact that these teas are most similar in terms

317    of processing conditions; they are subjected to either limited or no oxidation. In comparison,

318    oolong and black teas undergo oxidation during their manufacturing. This may be why these two

319    teas lie in two adjacent clusters that are far away from the cluster containing green, yellow and

320    white teas. However, further information is required to identify the exact source of the spectral

321    differences observed.

322

323    *Figure 5: Mahalanobis distances between the three different cluster groups obtained using PCA (dark-*

324    *blue), MDS (sky-blue), ISOMAP (light- green) and t-SNE (yellow).*

325    To assess further the separation of clusters with each method, the Mahalanobis distance between

326    the clusters was calculated. Figure 5 presents the Mahalanobis distance estimated for the three

327    major clusters identified in Figure 4. The three major cluster can be understood as the group of

328    minimally processed tea products available on the market (denoted the green group), the teas

329    subjected to oxidation (oxidised group) and those that have been subjected to microbial fermented

330    (fermented group). The x-axis in Figure 5 presents the pairwise groups used for estimating the

331    distance and the y-axis gives the respective Mahalanobis distance obtained from the different data

332    visualisation methods. It can be seen that the t-SNE (yellow) was superior to all other methods

333    followed by the ISOMAP (light green), and then PCA (dark blue) and MDS (sky blue) for

334    separating all three groups in the data-visualisation plots.

335    From a statistical perspective, a better visualisation of separate clusters corresponding to different

336    tea products with t-SNE could be due to its ability to capture the non-linearity present in the data

337    set and consideration of neighbourhood information. This supports the modelling of both distant

338    and nearby points (Maaten and Hinton, 2008). Often, in high dimensional space when the data lies

339    near, or in a non-linear manifold, linear methods like PCA and MDS fail to preserve the structure

340    of data in the lower dimension space. This is because with linear methods like PCA and MDS, the

341    aim is to keep the distant object far apart; no consideration is given to utilising the information

342    about the neighbouring data points (Maaten et al., 2009).

343    It can be seen in Figure 4(c) that ISOMAP provides a little insight on differences in the classes

344    belonging to black and oolong teas compared to what was achieved with PCA (Figure 4(a)) and

345    MDS (Figure 4(b)). However, ISOMAP was not able to provide a clear separation of the two teas

346    as was obtained with t-SNE. A reason for the poor performance of ISOMAP compared to t-SNE

347    could be due to its weakness in dealing with the holes and non-convex nature of the data manifold

348    in the higher dimension (Tenenbaum, 1998). Another important weakness of ISOMAP is its

349    topological instability, which leads to a short-circuiting problem in the neighbourhood graph and

350    results in its poor performance (Balasubramanian and Schwartz, 2002).

351    *3.3. Support vector machine classification*

352

*Figure 6: (a) Greyscale image constructed from the spectral plane extracted from the hypercube at 1424 nm, (b) Classification maps obtained from the application of the ECOC-SVM model. From left to right the samples can be understood as yellow (dark blue), oolong (light blue), green (cyan), black (light green), white (orange) and Pu-erh (yellow). (c) Histograms showing the proportion of pixels attributed to the different tea products for the classification maps in (b).*

The results from the application of the ECOC-SVM multi-class classification model are presented as classification maps in Figures 6 and 7. Figure 6(b) presents the classification maps of pure tea samples, from left to right, the samples can be understood as yellow, oolong, green, black, white and Pu-erh. For comparison, a greyscale image was also produced (Figure 6(a)) using the spectral plane corresponding to 1424 nm; this wavelength was selected merely to allow visualisation of the data hypercube. It can be seen from Figure 6(b) that all six teas were classified into their respective individual classes. However, there are some pixels that were misclassified; Figure 6(c) shows the proportion of pixels attributed to the different tea products for the classification maps in Figure 6(b). The misclassification was most dominant at the edges owing to signal from the circular sample container; such pixels (approximately 20%) were misclassified as Pu-erh. When these pixels were excluded, an overall accuracy of 97.41±0.16 % was obtained for cross-validated samples using 1200 iterations.

Apart from the edges, a reason for the misclassification between different teas can be attributed to their spectral similarity. When visualising the data with t-SNE (see Figure 4(d)), green, white and yellow tea were found to be lying near in the same cluster, and black and oolong were near to each other due to their spectral similarity. Hence, the classification map for the yellow tea (dark blue) has some misclassified pixels that have been attributed to either white (orange) or green tea (cyan). For black and oolong teas, it can be noted that there are some pixels in the classification map for

376  black tea (light green) that were misclassified as oolong (light blue class) and vice-versa. Another

377  possible reason for misclassification could arise from the purity of the tea; for example, a

378  minimally processed tea (e.g. white) may contain small amounts of oxidised product (e.g. black

379  tea).

380

381  *Figure 7: (a). Greyscale image at 1424 nm for the sample comprising oolong, black and Pu-erh tea, (b).*

382  *The classification map for the sample comprising oolong, black and Pu-erh tea, (c). Pie chart representing*

383  *the proportion of pixels belonging to a particular class for the classification map presented in (b), (d).*

384  *Greyscale image at 1424 nm for a sample containing a mixture of all teas, (e). The classification map for*

385  *a sample containing a mixture of all teas, and (f) Pie chart representing the proportion of pixels belonging*

386  *to a particular class for the classification map presented in (e).*

387  Figure 7 presents the classification maps for the HS images acquired for samples comprising

388  mixtures of teas. This analysis was performed to assess the feasibility of using the methodology

389  developed to classify different tea samples when more than one tea is present. Figure 7(a) presents

390  the spectral plane corresponding to 1424 nm for a sample containing oolong, black and Pu-erh teas

391  (not mixed) in roughly equal portions. These three teas were selected as there is an oxidation stage

392  in their manufacturing. The location of the different teas in Figure 7(a) can be identified with the

393  red markers. As can be seen from Figure 7(b), the model provided a clear classification of the three

394  teas into their respective classes. However, some misclassification can be seen at the interface

395  between different types of teas; individual pixels will detect the presence of more than one type of

396  tea at these locations. Furthermore, Figure 7(c) provides insight into the proportion of pixels

397  belonging to each class. It can be seen that the pie chart is mainly dominated by the proportion of

398  oolong, black and Pu-erh tea and contains a very small portion (<1 %) of pixels classified as green,

399    white and yellow.

400    The methodology developed was also tested for a mixture of all six tea samples. The result for

401    classification of the sample containing a mixture of all six types of tea is presented in Figure 7(e).

402    The classification map shown in figure 7(e) can be interpreted in conjunction with the pie chart

403    (Figure 7(f)) representing the proportion of pixels classified belonging to different classes. The pie

404    chart shows that the presence of all the classes can be detected with the classification model and

405    the portion of each type of tea ranged from 10 – 26%. However, it was not possible to validate the

406    classification result of the mixture image because it is not known if the sample was a homogenous

407    mixture of the six types of teas and hence, the exact composition of the upper surface of the sample

408    is unknown. In addition, there may be some misclassification of pixels that detect more than one

409    type of tea.

## 410   **4. Conclusions**

411    NIR HSI has been used to classify six different types of commercial tea samples. Before any data

412    modelling, the spectral imaging data from tea products should be pre-processed to reduce the

413    effects of light scattering arising from the inhomogeneous and uneven leaf surface. Four different

414    types of linear and non-linear dimensionality reduction methods were compared for visualisation

415    of imaging data. The non-linear method, t-SNE, gave better separation of the different tea products

416    than classical linear techniques such as PCA and MDS. This is because t-SNE uses information

417    from neighbouring data points in the high dimensional space to preserve the structure in the low

418    dimensional representation. It was possible to classify the tea according to product type using a

419    ECOC-SVM multi-class classification model constructed using the NIR HSI data. Therefore, NIR

420    HSI in conjunction with machine learning could be a potential tool for classification of different

421 types of tea products. The source of spectral differences is assumed to arise from the different

422 processing steps that are involved in the manufacture of various types of tea. However, there could

423 be other sources, e.g. geographical, that contribute to spectral differences and hence, this requires

424 further investigation.

## 5. Acknowledgments

430

## 6. References

432 Balasubramanian, M., Schwartz, E.L., 2002. The Isomap Algorithm and Topological

433 Stability. Science (80). 295, 7 LP –7.

434 Baldwin, E.A., Bai, J., Plotto, A., Dea, S., 2011. Electronic Noses and Tongues:

435 Applications for the Food and Pharmaceutical Industries. Sensors (Basel). 11(5), 4744-

436 4766. https://doi.org/10.3390/s110504744

437 Barnes, R.J., Dhanoa, M.S., Lister, S.J., 1989. Standard Normal Variate Transformation

438 and De-Trending of Near-Infrared Diffuse Reflectance Spectra. Appl. Spectrosc. 43,

439 772–777. https://doi.org/10.1366/0003702894202201

440 Bhattacharyya, N., Seth, S., Tudu, B., Tamuly, P., Jana, A., Ghosh, D., Bandyopadhyay,

441 R., Bhuyan, M., 2007. Monitoring of black tea fermentation process using electronic

442 nose. J. Food Eng. 80, 1146–1156.

443 https://doi.org/https://doi.org/10.1016/j.jfoodeng.2006.09.006

444 Bian, B.M., Skidmore, A.K., Schlerf, M., Fei, T., Liu, Y.F., Wang, T., 2010. Reflectance

445 spectroscopy of biochemical components as indicators of tea, Camellia Sinensis, quality

446 76, 1385–1392.

447 Bian, M., Skidmore, A.K., Schlerf, M., Wang, T., Liu, Y., Zeng, R., Fei, T., 2013.

448 Predicting foliar biochemistry of tea (Camellia sinensis) using reflectance spectra

449 measured at powder, leaf and canopy levels. ISPRS J. Photogramm. Remote Sens. 78,

450 148–156. https://doi.org/https://doi.org/10.1016/j.isprsjprs.2013.02.002

451 Chang, K, 2015. World tea production and trade: current and future development, Food

452 and Agriculture Organisation (FAO), United Nations (UN), Rome.

453 http://www.fao.org/3/a-i4480e.pdf (last accessed: 14 September 2017)

454 Chen, Q., Zhao, J., Cai, J., 2008. Identification of Tea Varieties Using Computer Vision.

455 Transactions of the ASABE. 51(2), 623-628. https://doi.org/10.13031/2013.24363

456 Chen, Q., Zhao, J., Fang, C.H., Wang, D., 2007. Feasibility study on identification of

457 green, black and Oolong teas using near-infrared reflectance spectroscopy based on

458 support vector machine (SVM). Spectrochim. Acta Part A Mol. Biomol. Spectrosc. 66,

459 568–574. https://doi.org/https://doi.org/10.1016/j.saa.2006.03.038

460       Chen, Q., Zhao, J., Huang, X., Zhang, H., Liu, M., 2006. Simultaneous determination of

461       total polyphenols and caffeine contents of green tea by near-infrared reflectance

462       spectroscopy. Microchem. J. 83, 42–47.

463       https://doi.org/http://dx.doi.org/10.1016/j.microc.2006.01.023

464       Chen, Q., Zhao, J., Lin, H., 2009. Study on discrimination of Roast green tea (Camellia

465       sinensis L.) according to geographical origin by FT-NIR spectroscopy and supervised

466       pattern recognition. Spectrochim. Acta Part A Mol. Biomol. Spectrosc. 72, 845–850.

467       https://doi.org/http://dx.doi.org/10.1016/j.saa.2008.12.002

468       Cox, T.F., Cox, M.A.A., 2000. Multidimensional scaling. CRC press.

469       Domínguez, I., Doménech-Carbó, A., 2015. Screening and authentication of tea varieties

470       based on microextraction-assisted voltammetry of microparticles. Sensors Actuators B

471       Chem. 210, 491–499. https://doi.org/https://doi.org/10.1016/j.snb.2015.01.009

472       Edelman, G.J., Gaston, E., van Leeuwen, T.G., Cullen, P.J., Aalders, M.C.G., 2012.

473       Hyperspectral imaging for non-contact analysis of forensic traces. Forensic Sci. Int. 223,

474       28–39. https://doi.org/https://doi.org/10.1016/j.forsciint.2012.09.012

475       Fu, X., Ying, Y., 2016. Food Safety Evaluation Based on Near Infrared Spectroscopy and

476       Imaging: A Review. Crit. Rev. Food Sci. Nutr. 56, 1913–1924.

477       https://doi.org/10.1080/10408398.2013.807418

478       García Allende, P.B., Anabitarte García, F., Conde Portilla, O.M., Mirapeix Serrano,

479       J.M., Madruga Saavedra, F.J., López Higuera, J.M., 2008. Support vector machines in

480  hyperspectral imaging spectroscopy with application to material identification. Proc.

481  SPIE 6966, Algorithms and Technologies for Multispectral, Hyperspectral, and

482  Ultraspectral Imagery XIV, 69661V. http://dx.doi.org/10.1117/12.770306

483  Garcia-Allende, P.B., Conde, O.M., Mirapeix, J., Cobo, A., Lopez-Higuera, J.M., 2008.

484  Quality control of industrial processes by combining a hyperspectral sensor and Fisher's

485  linear discriminant analysis. Sensors Actuators B Chem. 129, 977–984.

486  https://doi.org/https://doi.org/10.1016/j.snb.2007.09.036

487  Goetz, A.F., Vane, G., Solomon, J.E., Rock, B.N., 1985. Imaging spectrometry for Earth

488  remote sensing. Science 228, 1147–1153. https://doi.org/10.1126/science.228.4704.1147

489  He, Y., Li, X., Deng, X., 2007. Discrimination of varieties of tea using near infrared

490  spectroscopy by principal component analysis and BP model. J. Food Eng. 79, 1238–

491  1242. https://doi.org/https://doi.org/10.1016/j.jfoodeng.2006.04.042

492  Jing, J., Shi, Y., Zhang, Q., Wang, J., Ruan, J., 2017. Prediction of Chinese green tea

493  ranking by metabolite profiling using ultra-performance liquid chromatography–

494  quadrupole time-of-flight mass spectrometry (UPLC–Q-TOF/MS). Food Chem. 221,

495  311–316. https://doi.org/https://doi.org/10.1016/j.foodchem.2016.10.068

496  Kandpal, L.M., Lohumi, S., Kim, M.S., Kang, J.-S., Cho, B.-K., 2016. Near-infrared

497  hyperspectral imaging system coupled with multivariate methods to predict viability and

498  vigor in muskmelon seeds. Sensors Actuators B Chem. 229, 534–544.

499  https://doi.org/https://doi.org/10.1016/j.snb.2016.02.015

500       Kandpal, L.M., Tewari, J., Gopinathan, N., Boulas, P., Cho, B.-K., 2016. In-Process

501       Control Assay of Pharmaceutical Microtablets Using Hyperspectral Imaging Coupled

502       with Multivariate Analysis. Anal. Chem. 88, 11055–11061.

503       https://doi.org/10.1021/acs.analchem.6b02969

504       Kumar, A.S., Shanmugam, R., Nellaiappan, S., Thangaraj, R., 2016. Tea quality

505       assessment by analyzing key polyphenolic functional groups using flow injection analysis

506       coupled with a dual electrochemical detector. Sensors Actuators B Chem. 227, 352–361.

507       https://doi.org/https://doi.org/10.1016/j.snb.2015.12.072

508       Lee, M.-S., Hwang, Y.-S., Lee, J., Choung, M.-G., 2014. The characterization of caffeine

509       and nine individual catechins in the leaves of green tea (Camellia sinensis L.) by near-

510       infrared reflectance spectroscopy. Food Chem. 158, 351–357.

511       https://doi.org/10.1016/j.foodchem.2014.02.127

512       Li, J., Fu, B., Huo, D., Hou, C., Yang, M., Shen, C., Luo, H., Yang, P., 2017.

513       Discrimination of Chinese teas according to major amino acid composition by a

514       colorimetric {IDA} sensor. Sensors Actuators B Chem. 240, 770–778.

515       https://doi.org/http://dx.doi.org/10.1016/j.snb.2016.09.019

516       Lu, G., Fei, B., 2014. Medical hyperspectral imaging: a review. J. Biomed. Opt. 19,

517       10901. https://doi.org/10.1117/1.JBO.19.1.010901

518       Lv, H., Zhang, Y., Lin, Z., Liang, Y., 2013. Processing and chemical constituents of Pu-

519       erh tea: A review. Food Res. Int. 53, 608–618.

520       https://doi.org/https://doi.org/10.1016/j.foodres.2013.02.043

521        Maaten, L. van der, Hinton, G., 2008. Visualizing high-dimensional data using t-SNE. J.

522        Mach. Learn. Res. 9, 2579–2605.

523        Maaten, L. van der, Postma, E., Van den Herik, J., 2009. Dimensionality reduction: a

524        comparative review. Tilburg University Technical Report, TICC-TR 2009-005.

525        Mahalanobis, P.C., 1936. On the generalised distance in statistics, in: Proceedings

526        National Institute of Science, India. pp. 49–55.

527        Mishra, P., Asaari, M.S.M., Herrero-Langreo, A., Lohumi, S., Diezma, B., Scheunders,

528        P., 2017. Close range hyperspectral imaging of plants: A review. Biosyst. Eng. 164, 49–

529        67. https://doi.org/https://doi.org/10.1016/j.biosystemseng.2017.09.009

530        Mishra, P., Cordella, C.B.Y., Rutledge, D.N., Barreiro, P., Roger, J.M., Diezma, B.,

531        2016. Application of independent components analysis with the JADE algorithm and NIR

532        hyperspectral imaging for revealing food adulteration. J. Food Eng. 168, 7–15.

533        Mountrakis, G., Im, J., Ogole, C., 2011. Support vector machines in remote sensing: A

534        review. ISPRS J. Photogramm. Remote Sens. 66, 247–259.

535        https://doi.org/https://doi.org/10.1016/j.isprsjprs.2010.11.001

536        Nansen, C., Singh, K., Mian, A., Allison, B.J., Simmons, C.W., 2016. Using

537        hyperspectral imaging to characterize consistency of coffee brands and their respective

538        roasting classes. J. Food Eng. 190, 34–39.

539        https://doi.org/https://doi.org/10.1016/j.jfoodeng.2016.06.010

540        Nieh, C.-H., Hsieh, B.-C., Chen, P.-C., Hsiao, H.-Y., Cheng, T.-J., Chen, R.L.C., 2009.

541        Potentiometric flow-injection estimation of tea fermentation degree. Sensors Actuators B

542        Chem. 136, 541–545. https://doi.org/https://doi.org/10.1016/j.snb.2008.09.024

543        Ozturk, B., Seyhan, F., Ozdemir, I.S., Karadeniz, B., Bahar, B., Ertas, E., Ilgaz, S., 2016.

544        Change of enzyme activity and quality during the processing of Turkish green tea. LWT -

545        Food Sci. Technol. 65, 318–324. https://doi.org/https://doi.org/10.1016/j.lwt.2015.07.068

546        Panigrahi, N., Bhol, C.S., Das, B.S., 2016. Rapid assessment of black tea quality using

547        diffuse reflectance spectroscopy. J. Food Eng. 190, 101–108.

548        https://doi.org/https://doi.org/10.1016/j.jfoodeng.2016.06.020

549        Plaza, A., Benediktsson, J.A., Boardman, J.W., Brazile, J., Bruzzone, L., Camps-Valls,

550        G., Chanussot, J., Fauvel, M., Gamba, P., Gualtieri, A., Marconcini, M., Tilton, J.C.,

551        Trianni, G., 2009. Recent advances in techniques for hyperspectral image processing.

552        Remote Sens. Environ. 113, S110–S122.

553        https://doi.org/https://doi.org/10.1016/j.rse.2007.07.028

554        Pu, Y.-Y., Feng, Y.-Z., Sun, D.-W., 2015. Recent Progress of Hyperspectral Imaging on

555        Quality and Safety Inspection of Fruits and Vegetables: A Review. Compr. Rev. Food

556        Sci. Food Saf. 14, 176–188. https://doi.org/10.1111/1541-4337.12123

557        Qu, J.-H., Liu, D., Cheng, J.-H., Sun, D.-W., Ma, J., Pu, H., Zeng, X.-A., 2015.

558        Applications of Near-infrared Spectroscopy in Food Safety Evaluation and Control: A

559        Review of Recent Research Advances. Crit. Rev. Food Sci. Nutr. 55, 1939–1954.

560        https://doi.org/10.1080/10408398.2013.871693

561    Ruan, J., Haerdter, R., Gerendás, J., 2010. Impact of nitrogen supply on carbon/nitrogen

562    allocation: a case study on amino acids and catechins in green tea [Camellia sinensis (L.)

563    O. Kuntze] plants*. Plant Biol. 12, 724–734. https://doi.org/10.1111/j.1438-

564    8677.2009.00288.x

565    Sabale, S.P., Jadhav, C.R., 2015. Hyperspectral Image Classification Methods in Remote

566    Sensing - A Review, in: 2015 International Conference on Computing Communication

567    Control and Automation. pp. 679–683. https://doi.org/10.1109/ICCUBEA.2015.139

568    Sang, S., 2016. Tea: Chemistry and Processing BT  - Encyclopedia of Food and Health.

569    Academic Press, Oxford, pp. 268–272. https://doi.org/https://doi.org/10.1016/B978-0-12-

570    384947-2.00685-1

571    Savitzky, A., Golay, M.J.E., 1964. Smoothing and Differentiation of Data by Simplified

572    Least Squares Procedures. Anal. Chem. 36, 1627–1639.

573    https://doi.org/10.1021/ac60214a047

574    Sharma, P., Ghosh, A., Tudu, B., Sabhapondit, S., Baruah, B.D., Tamuly, P.,

575    Bhattacharyya, N., Bandyopadhyay, R., 2015. Monitoring the fermentation process of

576    black tea using QCM sensor based electronic nose. Sensors Actuators B Chem. 219, 146–

577    157. https://doi.org/https://doi.org/10.1016/j.snb.2015.05.013

578    Shrestha, S., Knapič, M., Žibrat, U., Deleuran, L.C., Gislum, R., 2016. Single seed near-

579    infrared hyperspectral imaging in determining tomato (Solanum lycopersicum L.) seed

580    quality in association with multivariate data analysis. Sensors Actuators B Chem. 237,

581    1027–1034. https://doi.org/https://doi.org/10.1016/j.snb.2016.08.170

582  Tan, J., Dai, W., Lu, M., Lv, H., Guo, L., Zhang, Y., Zhu, Y., Peng, Q., Lin, Z., 2016.

583  Study of the dynamic changes in the non-volatile chemical constituents of black tea

584  during fermentation processing by a non-targeted metabolomics approach. Food Res. Int.

585  79, 106–113. https://doi.org/https://doi.org/10.1016/j.foodres.2015.11.018

586  Tan, S.-M., Luo, R.-M., Zhou, Y.-P., Gong, H., Tan, Z., 2012. Rapid and non-destructive

587  discrimination of tea varieties by near infrared diffuse reflection spectroscopy coupled

588  with classification and regression trees. African J. Biotechnol. 11, 2303–2312.

589  Tenenbaum, J.B., 1998. Mapping a Manifold of Perceptual Observations, in: Jordan,

590  M.I., Kearns, M.J., Solla, S.A. (Eds.), Advances in Neural Information Processing

591  Systems 10. MIT Press, pp. 682–688.

592  Vapnik, V.N., Vapnik, V., 1998. Statistical learning theory. Wiley New York.

593  Wang, S., Yang, X., Zhang, Y., Phillips, P., Yang, J., Yuan, T.-F., 2015. Identification of

594  Green, Oolong and Black Teas in China via Wavelet Packet Entropy and Fuzzy Support

595  Vector Machine. Entropy. 17, 6663-6682. https://doi.org/10.3390/e17106663

596  Wang, X., Huang, J., Fan, W., Lu, H., 2015. Identification of green tea varieties and fast

597  quantification of total polyphenols by near-infrared spectroscopy and ultraviolet-visible

598  spectroscopy with chemometric algorithms. Anal. Methods 7, 787–792.

599  https://doi.org/10.1039/C4AY02106A

600  Wold, S., Esbensen, K., Geladi, P., 1987. Principal component analysis. Chemom. Intell.

601  Lab. Syst. 2, 37–52. https://doi.org/https://doi.org/10.1016/0169-7439(87)80084-9

602       Xie, C., Li, X., Shao, Y., He, Y., 2015. Color Measurement of Tea Leaves at Different

603       Drying Periods Using Hyperspectral Imaging Technique. PLoS One 9, 1–15.

604       https://doi.org/10.1371/journal.pone.0113422

605       Yaroshenko, I., Kirsanov, D., Kartsova, L., Bhattacharyya, N., Sarkar, S., Legin, A.,

606       2014. On the application of simple matrix methods for electronic tongue data processing:

607       Case study with black tea samples. Sensors Actuators B Chem. 191, 67–74.

608       https://doi.org/https://doi.org/10.1016/j.snb.2013.09.093

609       Zhao, J., Chen, Q., Cai, J., Ouyang, Q., 2009. Automated tea quality classification by

610       hyperspectral imaging. Appl. Opt. 48, 3557–3564. https://doi.org/10.1364/AO.48.003557

611       Zhu, H., Ye, Y., He, H., Dong, C., 2017. Evaluation of green tea sensory quality via

612       process characteristics and image information. Food Bioprod. Process. 102, 116–122.

613       https://doi.org/https://doi.org/10.1016/j.fbp.2016.12.004

*Figure 2: Processing steps for different tea products starting from fresh green tea leaves to final products.*

621

622 *Figure 2: Illustrative diagram for the hyperspectral imaging setup used to acquire the images of*
623 *tea samples.*

624



625

626 *(a)*            *(b)*            *(c)*

627 *Figure 3: Absorbance spectra of pure tea samples of yellow, oolong, green, black, white and Pu-*
628 *erh. (a). Mean absorbance spectra (n = 200). (b) Mean spectra after pre-processing (SNV and*
629 *Savitzky- Golay smoothing), and (c) standard deviation of the absorbance spectra and spectra*
630 *after pre- processing. The vertical green lines denote the positions of the main peaks.*

631

632

633

*Figure 4: 2-Dimensional scatter plots for visualising high dimensional tea data. (a). Principal Component Analysis (PCA), (b). Multidimensional Scaling (MDS), (c). Isometric Mapping (ISOMAP), and (d). t-distributed Stochastic Neighbour Embedding (t-SNE). In all the plots, the first dimension is represented in the x-axis and the second in the y-axis, and the six tea products are coloured as follows: Pu-erh (pink), black (sky blue), oolong (yellow), green (green), white (blue) and yellow (red).*

640

641

642  *Figure 5: Mahalanobis distances between the three different cluster groups obtained using PCA*

643  *(dark-blue), MDS (sky-blue), ISOMAP (light- green) and t-SNE (yellow).*

644

*Figure 6: (a) Greyscale image constructed from the spectral plane extracted from the hypercube at 1424 nm, (b) Classification maps obtained from the application of the ECOC-SVM model. From left to right the samples can be understood as yellow (dark blue), oolong (light blue), green (cyan), black (light green), white (orange) and Pu-erh (yellow). (c) Histograms showing the proportion of pixels attributed to the different tea products for the classification maps in (b).*
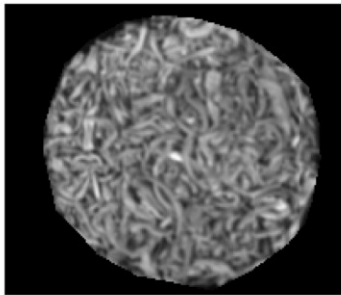
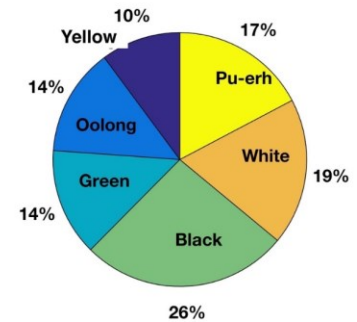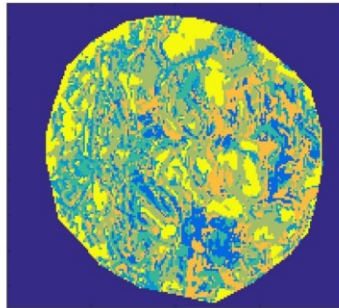Figure 7: (a). Greyscale image at 1424 nm for the sample comprising oolong, black and Pu-erh tea, (b). The classification map for the sample comprising oolong, black and Pu-erh tea, (c). Pie chart representing the proportion of pixels belonging to a particular class for the classification map presented in (b), (d). Greyscale image at 1424 nm for a sample containing a mixture of all teas, (e). The classification map for a sample containing a mixture of all teas, and (f) Pie chart representing the proportion of pixels belonging to a particular class for the classification map presented in (e).