This is a repository copy of *Rapid identification of Lactobacillus species using near infrared spectral features of bacterial colonies*.

White Rose Research Online URL for this paper:
https://eprints.whiterose.ac.uk/179320/

Version: Accepted Version

**Article:**

1    # Rapid identification of *Lactobacillus* strains using near-infrared

2    # spectroscopy and chemometrics

3    Jiyong Shi[a], Xuetao Hu[a], Fang Zhang[a], Xiaobo Zou[a, *], Mel Holmes[b], Zhihua Li[a],

4    Yiwei Xu[a], Wen Zhang[a], Xiaowei Huang[b]

5    [a] School of Food and Biological Engineering, Jiangsu University, Zhenjiang 212013,

6    China

7    [b] School of Food Science and Nutrition, the University of Leeds, Leeds LS2 9JT,

8    United Kingdom

9    *Corresponding author. Tel: +86 511 88780085; Fax: +86 511 88780201

10   Email address: zou_xiaobo@ujs.edu.cn (xiaobo Zou)

11   Highlights

12   NIR spectroscopy was applied for identification of *Lactobacillus* strains according to

13   the discrepancy of NIR acquired from bacterial colonies.

14   PCA and HCA were investigated for general discriminant analysis of *Lactobacillus*.

15   UVE and GA were performed for selection of characteristic wavenumbers.

16   Sensitivity in identification of *Lactobacillus* strains was excellent (92.857%) using

17   UVE-GA-LS-SVM with 10 wavenumbers.

18   IR spectra of bacterial colonies and cells were researched and compared in order to

19   demonstrate the feasibility of using NIR spectra acquired from bacterial colonies for

20   identification of *Lactobacillus*.

21

22    Abstract:

23        *Lactobacillus* (*L.*) plays an important role in food fermentation, while the

24    presence of some particular *Lactobacillus* species may decrease the quality of

25    fermented food products. In this study, near-infrared (NIR) spectral features of

26    *Lactobacillus* species were extracted and the feasibility in rapid identification of

27    *Lactobacillus* based on these NIR spectral features was investigated. Bacterial

28    colonies of four *Lactobacillus* species (*L. breris*, *L. casei*, *L. fermentum*, *L. reuteri*)

29    were cultured using spread-plate technique with MRS agar medium. Raw NIR

30    spectral data of bacterial colonies were acquired in the wavelength range of

31    4,000-10,000 cm$^{-1}$. After pre-processing, uninformative variables elimination (UVE)

32    and genetic algorithm (GA) were used to select the characteristic wavelengths

33    correlated with the four species of *Lactobacillus*. The NIR data corresponding to the

34    selected characteristic wavelengths were employed to build identification models for

35    discriminating the four species of *Lactobacillus* using Least Squares Support Vector

36    Machine (LS-SVM). The recognition rates of calibration set and prediction set using

37    the optimal identification model were 100% and 92.857%, respectively. In order to

38    explain scientifically the good results of NIR, mid infrared (MIR) spectra of bacterial

39    cells were collected and analyzed. Analytical results of MIR indicated that (1)

40    significant differences were observed in MIR spectral data collected from the four

41    species of *Lactobacillus;* (2) the selected NIR wavenumbers were quite correlated

42    with the MIR wavelengths that could reflect changes of components and structure

43    among four *Lactobacillus* cells. This explained why the four species of *Lactobacillus*

44    were reasonably well identified based on its NIR data. It is concluded that NIR

45    spectroscopy combining with chemometrics methods proposed in this paper could be

46    applied for rapid discrimination of *Lactobacillus*.

47    Keywords: NIR spectroscopy; *Lactobacillus*; Uninformative variables elimination:

48    Genetic algorithm; Least Squares support vector machine

49

## 1. Introduction

The lactic acid bacteria (LAB) are rod-shaped bacilli or cocci characterized by an increased tolerance to a lower pH range[1]. Various genera of lactic acid bacteria used as probiotics can improve the intestinal immune status, and maintain microbial balance during gastrointestinal disturbances[2]. *Lactobacillus* (*L.*) is one of the most important genera in LAB fermenting glucose primarily to lactic acid, $CO_2$ and ethanol[3]. *Lactobacillus* are widely used in food fermentation of yogurt[4], cheese[5], soybean meal[6], etc.[7, 8]. However, some *Lactobacillus* species, such as *L. brevis*, have been recognized as spoilage bacteria in food processes such as beer fermentation[9]. Uncontrolled growth of *Lactobacillus* in beer fermentation process may decrease the organoleptic quality of the beer (turbidity, sediment, acidification, off-flavor and ropiness), and even affect its safety[9, 10]. Therefore, it is necessary to identify the specie of *Lactobacillus* in food production.

There is differentiation in chemical and physical properties between bacterial strains so that they can be differentiated from one another. Traditional approaches for identifying *Lactobacillus* are phenotypic method (*5, 11, 12*) and molecular biological method (*13-15*). In the phenotypic method, a lot of phenotypic bacteria characteristics are used to identify an isolate of *Lactobacillus*. These phenotypic bacteria characteristics include morphological features[16] (shape, size, color, dimensions, form, etc.), physiological features[5] (modes of fermentation, acid/alkaline/salt tolerance, content of lactic acid dehydrogenases, etc.), and biochemical features[17] (the type of interbridging in the peptidoglycan, type of fatty acids and proteins, etc.). In general, about 17 phenotypic tests are required to acquire these phenotypic bacteria characteristics, which makes the phenotypic method is time-consuming, tedious, and involve numerical preparation procedures[18]. Researchers also pointed out that the reliability of these tests has been questioned[16]. The exact identifications of these closely related species were not reliable; some were doubtful or unacceptable and some strains were misidentified with a good identification level[19]. In the molecular biological method, different molecular features based on molecular characterization

79  techniques (genotyping, multilocus sequence typing (MLST)(*20*), pulsed-field gel

80  electrophoresis (PFGE)(*21*) , and ribotyping(*22*)) are used to identify an isolate of

81  *Lactobacillus(19)*. The molecular biological method has been a powerful tool that

82  have helped microbiologists to detect the smallest variations within microbial species

83  even within individual strains. However, the molecular biological method requires

84  advanced instruments and very well trained hands for elaborated sample preparation

85  (*14, 23*). Therefore, there is a need to investigate a rapid and simple technology for

86  identification of bacterial strains.

87  Usually, optical methods, such as mid infrared (MIR) spectroscopy and near

88  infrared (NIR) spectroscopy, can obtain bacteria characteristics more quickly than

89  phenotypic method and molecular biological method. Optical methods have been

90  employed to rapidly identify bacterial strains (*24-27*). Recently, MIR and NIR were

91  used to record the spectral features of different bacterial strains, and bacterial strains

92  were identified based on these spectral features due to the fact that these spectral

93  features could reflect differences in chemical components (proteins, fatty acids,

94  nucleic acids, etc.)(*28-30*) between bacterial strains. Various researches have shown

95  that MIR can be used to differentiate and identify a number of microorganism at

96  different taxonomic levels, such as, *Lactobacilli* strains(*31*), *Filamentous Fungi*

97  strains(*26*), *Brettanomyces bruxellensis* strains(*32*) et.al. However, it is difficult to

98  acquire MIR spectra of bacteria directly owing to bacterial cell which could not be

99  seen by naked eye. In their researches, the following procedures were applied for

100  preparation of bacterial sample in order to obtain the corresponding MIR spectra of

101  strains: (1) strain was placed into MRS broths and incubated; (2) bacterial broth was

102  centrifuged and supernatant discarded; (3) the bacterial pellet was dried under

103  moderate vacuum to obtain a transparent bacterial film before spectral acquisition.

104  Due to the complicated and time-consuming bacterial preparation which are

105  indispensable before acquisition of MIR spectra, it is no longer a simple and fast

106  method to identify strains using MIR spectroscopy. NIR spectroscopy technique, as an

107  alternative rapid detection method, has been investigated the feasibility in

108  identification and classification of *Escherichia coli* strains(*30*), *Pseudomonas*

109    *aeruginosa* strains(*33*), *Bacillus amyloliquifaciens* strains(*34*), *Bacillus cereus*

110    strains(*34*), *Listeria innocua* strains(*30)* et.al. In order to overcome the difficulties of

111    spectral acquisition using the NIR spectroscopy on account of strains with small size,

112    the NIR spectra were acquired from bacterial suspension after resulting pellets

113    re-suspended into a series of serial dilutions. The resulting pellets were obtained from

114    the centrifugation of bacterial broths after inoculation and culture. Although the

115    sample preparation for acquisition of NIR spectra is simpler than the sample

116    preparation for acquisition of MIR spectra. It is difficult to identify all species at the

117    same time when more than one specie simultaneously exist in the bacterial suspension.

118    Therefore, it would restrict the application of NIR in rapid identification of strains.

119    　　It is noted that the colony is enriched with a large number of bacterial cells

120    owing to self-replication of single strain and the classification of different strains

121    maybe performed by the unique NIR spectra of each strain. When the bacterial

122    colonies at different species are cultured on agar plate simultaneously, NIR spectra of

123    each colony are obtained. Therefore, rapid and simultaneous classification of multiple

124    strains could be performed by the NIR spectra. To validate the feasibility of this idea,

125    the following researches were investigated: (1) NIR spectra of four *Lactobacillus*

126    were acquired from bacterial colonies by NIR spectroscopy; (2) methodology for

127    rapid identification of *Lactobacillus* was developed by NIR spectral data and

128    chemometrics; (3) MIR spectra were obtained to prove the validation of identification

129    using the NIR spectra acquired from bacterial colony.

130    ## 2. Materials and methods

131    2.1 Strains and preparation of samples

132    　　Four species of *Lactobacillus* (*L. breris*, *L. casei*, *L. fermentum*, *L. reuteri*) from

133    Fermentation Laboratory of the school of Food and Biological Engineering, Jiangsu

134    University were used in this study. Four *Lactobacillus* were stored in glycerol at

135    -80°C prior to use. *Lactobacillus* grew on a tube containing 25 mL

136    deMan-Rogosa-Sharpe (MRS) broth and incubated at 37°C for 24h. Each broth was

137    diluted at a ratio of 1:9 (broth: distilled water). Following repeated dilution, 1 ml of

138   the each dilution was extracted and smeared with a sterilized spreader onto MRS agar

139   plates, which was incubated at 37°C for 36h. Three replicates from each *Lactobacillus*

140   were cultivated and prepared in dependent assays to produce three independent

141   sample. Therefore, 21 plates for each *Lactobacillus* were obtained. Bacterial colonies

142   in plates were used to collect NIR spectral data for further analysis.
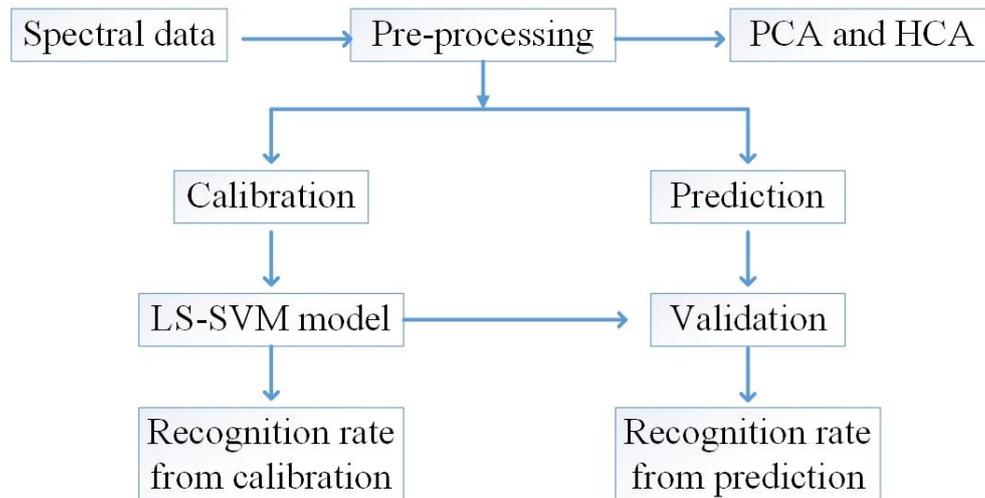
143   2.2 NIR spectral acquisition

144       The Antaris$^{TM}$ II near-infrared spectrophotometer (Thermo Electron Corporation,

145   USA) was employed for spectral acquisition over the 10,000 cm$^{-1}$ to 4000 cm$^{-1}$ at

146   intervals of 3.856 cm$^{-1}$ which resulted in 1557 variables. The plates, containing

147   individual *Lactobacillus* colonies, was placed on sample stage of NIR spectroscopy

148   for acquisition of bacterial spectra. Five spectra were collected from five

149   discontinuous region in individual cultured plate with the diffuse reflectance

150   integrating sphere mode. The average of five spectra from same plate was computed

151   and used in further data processing. Finally, 84 NIR spectra (21 spectra for each

152   *Lactobacillus*) from 84 plates respectively can be considered independent samples and

153   used for identification of *Lactobacillus*.

154   2.3 Chemometrics methods

155       Due to the NIR spectra significantly influenced by non-linearities and baseline

156   shift introduced by light scatter, suitable pre-processing should be applied to eliminate

157   these effects largely. For spectral data compression and information extraction, PCA

158   were performed to transform the spectral data after pre-processing into new

159   uncorrelated variable called principal components (PCs). HCA was applied with the

160   PCs as input variables to investigate similarities between different bacterial samples

161   and reveal if there are natural clustering in bacterial samples. Employing the full

162   spectral data does not always yield optimal results as the full variables may include

163   variables not correlated with colony features, partially correlated with colony features

164   and highly correlated with colony features. Therefore, UVE combined with GA called

165   UVE-GA is employed for characteristic variables selection in which UVE is

166   employed for elimination of uninformative variables and GA is for selection of

167   characteristic variables. LS-SVM was implemented in this study for classification

168    purposes. After the LS-SVM models were developed, they were applied to classify

169    bacterial samples. The data analysis mentioned above were performed by MATLAB

170    2010a after NIR spectral collection and steps of data analysis were shown in Fig.1.

171



172        Fig.1 Flowchart of data analysis in identification of *Lactobacillus*

173    2.3.1 Pre-processing methods

174        All pre-processing technique have the goal of reducing the un-modeled

175    variability noise in data in order to enhance the feature in the spectra. However, there

176    is always the danger of applying the wrong type or applying a severe processing that

177    will remove the valuable information(*35*). In this study, standard normal variate

178    (SNVT), multiplicative scatter correction (MSC), first derivative (1D) and second

179    derivative (2D) were performed to correct for light scattering, modify the additive and

180    multiplicative effects, and remove the influence of any baseline variation(*36*).

181    2.3.2 PCA and HCA

182        The PCA is a technique that is used for spectral data compression and

183    information extraction by the reduction of the data dimensionality and generation of

184    new uncorrelated variables called PCs(*37*). The PCs that account for a large

185    percentage of total variance reveal most information of *Lactobacillus* (*30*). HCA is a

186    unsupervised pattern recognition technique that can be used for clustering of strains

187    based on the similarities between spectra (*37*). The HCA was performed with the PCs

188    as input data using Ward's clustering algorithm and the Squared Euclidean Distance

189    Measure to generate a dendrogram(*38, 39*).

190 2.3.3 Variable selection methods (UVE, GA)

191     Uninformative variables elimination (UVE) is employed to eliminate

192 uninformative variables that clearly have no information about *Lactobacillus* in this

193 study. The artificial random variables are added to the calibration data as a reference

194 so that those spectral variables which play a less important role in model than the

195 random variables are eliminated(*40*). Compared to other methods, it keeps a large

196 number of the partially relevant variables for finally modeling. Therefore, GA was

197 employed to reduce the number of variables in order to achieve a simpler, more robust

198 model. GA is a popular heuristic optimization technique that employs a probabilistic,

199 non-local search process. Many studies have demonstrated the importance of GA for

200 characteristic wavelength selection(*41*).

201 2.3.4 Pattern recognition method (LS-SVM)

202     Before building classification models, the spectra data were randomly assigned

203 to calibration and test set according to the proportion of 2:1. For supervised

204 classification models, least squares support vector machine (LS-SVM) was

205 implemented. By applying LS-SVM classifiers, the empirical classification error can

206 be minimized while ensuring maximization of the interclass geometric boundary(*37*).

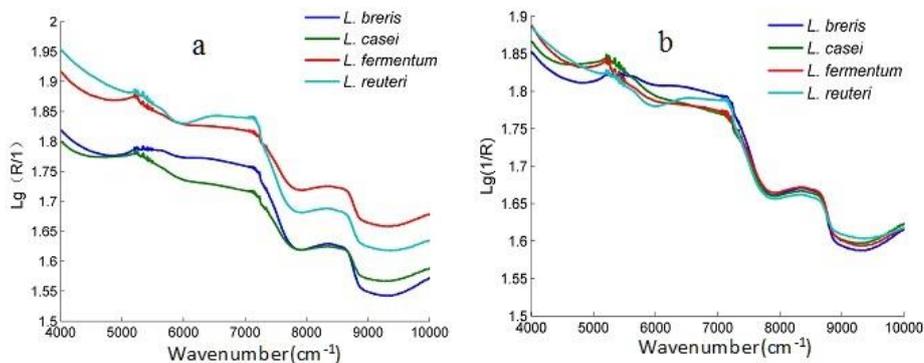207 2.3 Development of identification models

208     The full-spectrum LS-SVM model using calibration set was obtained based on

209 full spectrum and LS-SVM. The UVE-LS-SVM model was obtained based on the

210 characteristic wavelengths selected by UVE. The characteristic wavelengths selected

211 by UVE and GA were applied to build UVE-GA-LS-SVM model. And then models

212 are validated by test sets. The performance of model are evaluated via recognition

213 rates of calibration and test sets(*30*).

214 **3. Result and discussion**

215 3.1 NIR spectra investigation

216     The spectra in the range of 4000-10000 cm$^{-1}$ illustrated that the absorption

217 decrease with the wavenumber increasing expect some spectral peaks (Fig.1). As can

218 be seen in Fig.1 (a), average NIR spectra of each *Lactobacillus* has no significant

219 discrepancy instead of the intensity. Three obvious absorption peaks around 5350

220　cm$^{-1}$, 7150 cm$^{-1}$ and 8650 cm$^{-1}$ are all based on molecular overtone and combination

221　vibrations. MSC, 1D, 2D, SNVT, MSC+1D, MSC+2D, SNV+1D and SNV+2D were

222　applied for the multivariate data to eliminate spectrum interference factors due to the

223　influence of noise or the diffuse reflectance. LS-SVM based on raw spectral data and

224　the spectral data after eight kinds of pre-processing methods was employed to select

225　most suitable pre-processing method according to the recognition rates showing in

226　Table 1. It was found that the MSC performed relatively better than the other noted

227　methods for separation of the *Lactobacillus* (Table 1). Therefore, MSC was selected

228　to be employed for further process in this study. The average spectra of each

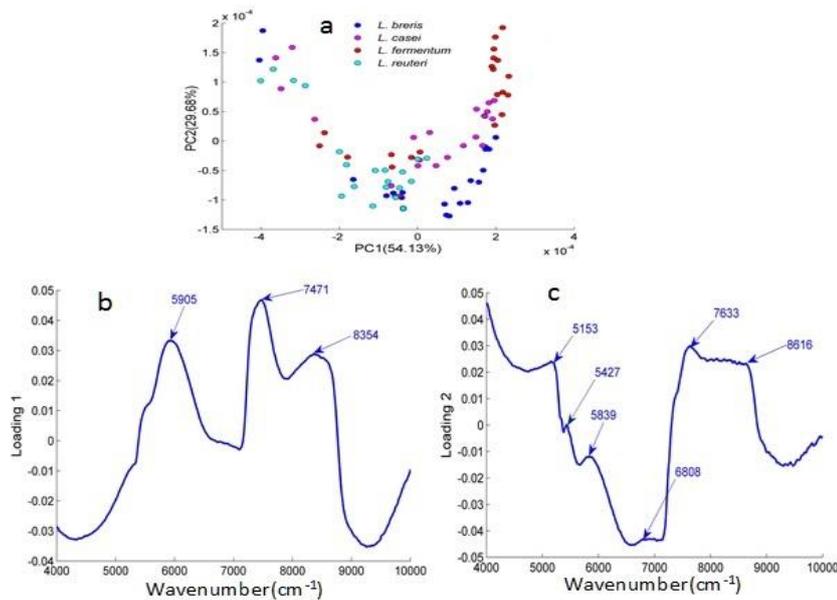229　*Lactobacillus* after MSC was shown in Fig.1 (b).

230



231　Fig.1. the raw average NIR spectra of four strains (a) and the NIR spectra after MSC (b)

232　3.2 Clustering analysis using PCA and HCA

233　　The NIR spectral data of four *Lactobacillus* (81 objects × 1557 variables) after

234　MSC were submitted to PCA to extract the effective information. New uncorrelated

235　variables were created, and the scores plot of PC1 versus PC2 display clustering for

236　*Lactobacillus* (Fig.2a). The first two principal components together accounted for

237　approximately 83.81% of the total variance (PC1: 54.13% and PC2: 29.68%). Based

238　on Fig.2a, a number of spectra for each *Lactobacillus* are overlapped. This weak

239　separation was obtained and maybe owing to the first two PCs that contained

240　insufficient information for classification. Therefore, more PCs and more

241　chemometrics, such as variable selection, supervised pattern recognition, should be

242　attempted to identify *Lactobacillus*. The loading plots of the first two PCs were

243　investigated to show wavenumbers making great contribution to the variation in the

244    data set. In Fig.2 (b), the peaks at 5905cm$^{-1}$, 7471cm$^{-1}$, 8354cm$^{-1}$ influencing the PC1

245    significantly could be the important wavenumbers that have high correlation with the

246    features of each *Lactobacillus*. Furthermore, the PC2 also was responsible for the

247    separation of *L. breris* and *L. fermentum*. The corresponding loading plot of PC2

248    shown in Fig.2 (c) manifest that 5153cm$^{-1}$, 5427cm$^{-1}$, 5839cm$^{-1}$, 6808cm$^{-1}$, 7633cm$^{-1}$,

249    8616cm$^{-1}$, associated with O-H, N-H structure, are the critical wavenumbers for such

250    separation.



251

252    Fig.2. PCA results: (a) score plot of PC1 versus PC2, (b) the loading plot for PC1, and (c) the loading plot for PC2

253        Fig.3 shows the dendrogram obtained after carrying out the HCA. The first ten

254    PCs, which account for 90% of total variance in dataset were used as the input data.

255    Ward's algorithm and squared Euclidian distance was selected to investigate the

256    dissimilarities between the spectra of four strains. In Fig.3, the left vertical axis of

257    dendrogram depicts the names of 40 bacterial samples. The second column on the left

258    is the number from 1 to 40, in which from 1 to 10 represent *L. breris*, 11 to 20

259    represent *L. casei*, 21 to 30 represent *L. fermentum*, and 31 to 40 represent *L. reuteri*.

260    The upper horizontal axis represents the distance between two bacterial samples or

261    two cultures. The magnitude of this distance depends on the number of spectra in a

262    cluster and the similarities between them. Two major clusters were illustrated in

263    dendrogram. The first cluster (two clusters) only included *L. breris* and the second

264 cluster was made up of two well distinguished subclusters, the first subcluster with *L.*

265 *fermentum* (10 strain), the second subcluster with *L. reuteri* (10 strains) and *L. casei*

266 (10 strains). The second subcluster is divided into two clusters, one is *L. reuteri*, and

267 the other one is *L. casei*. Only single *L. reuteri* was wrongly assigned to the cluster of

268 *L. casei*. The dendrogram in this paper does not completely follow the classical

269 scheme (homofermentative and heterofermentative *Lactobacillus*) using FT-IR

270 spectroscopy(*29*). In addition, the dendrogram provided in this study was also not

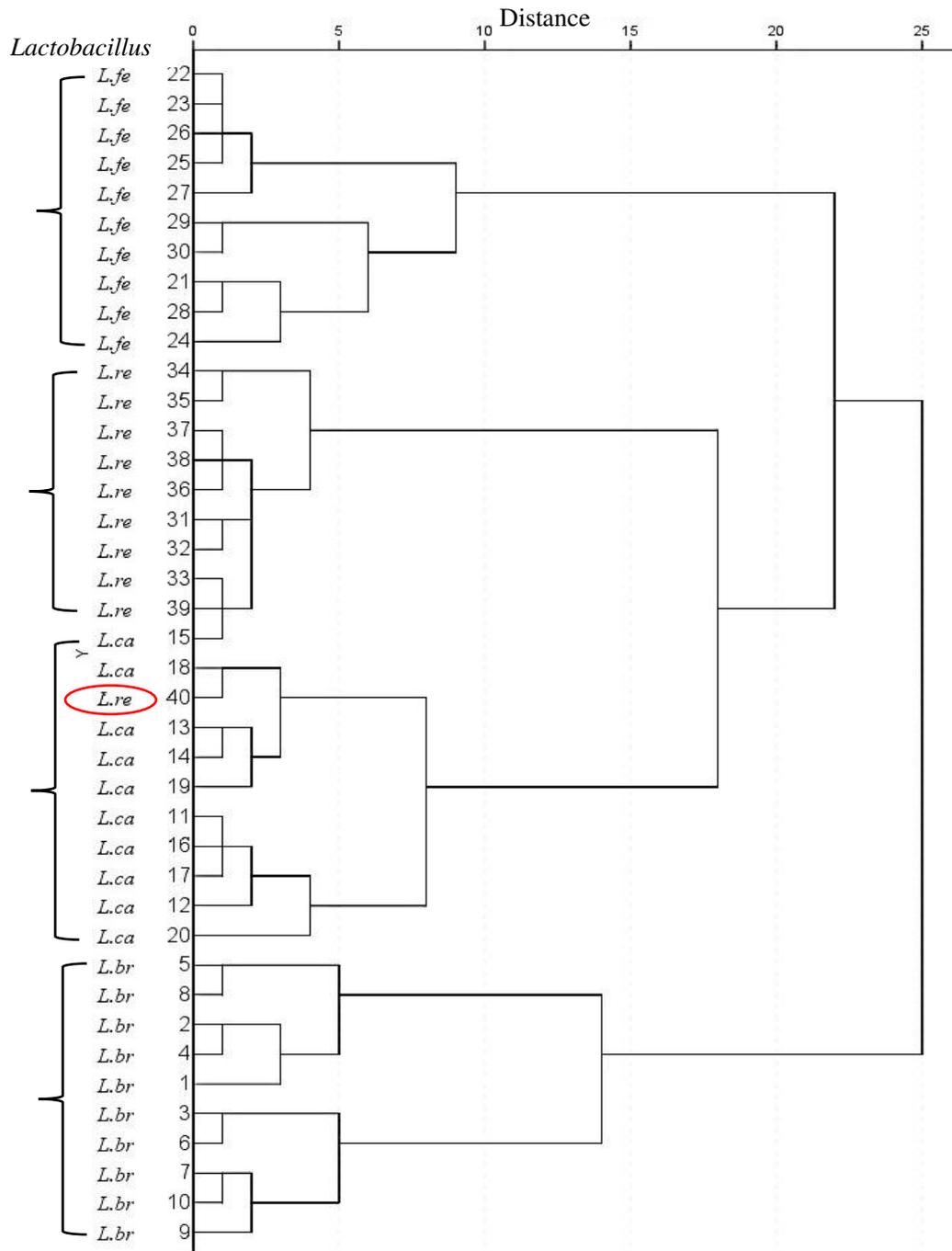271 similar with phylogenetic tree of *Lactobacillus* (*42, 43*),

Fig.3 Dendrogram from hierarchical cluster analysis of four *Lactobacillus*

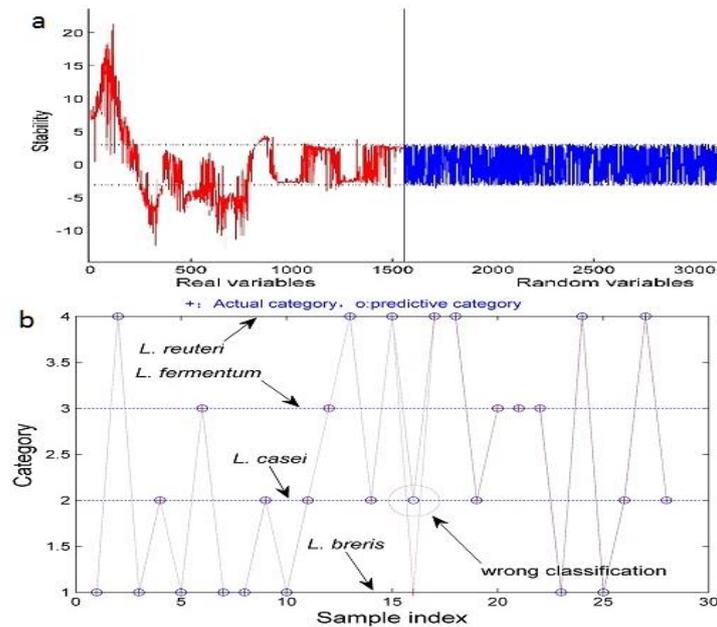3.3 Development of identification models using UVE, GA, and LS-SVM

In this study, 56 samples were used for calibration and 28 samples were used for prediction data. Classification results based on full spectrum data after eight different pre-processing were illustrated in Tab.1. As indicated in Tab.1, the highest accuracy of LS-SVM model was based on MSC with recognition rate of 85.714% for test set and 98.214% for the calibration set. The second high recognition rate for prediction

280   set is 82.143% when SNV is used. In the best model, four bacterial samples were

281   mistakenly classified. One *L. fermentum* was mistakenly predicted as *L. casei*, one *L.*

282   *breris* was mistakenly predicted as *L. fermentum*, one *L. fermentum* was mistaken for

283   *L. reuteri*, and *one L. casei* was mistaken for *L. brevis*.

284       Owing to the high-dimensional data containing highly correlated variables, UVE

285   and GA were employed in this research for selecting characteristic wavenumbers

286   reflecting the features of *Lactobacillus*. It could be advantageous to use only few

287   variables for accurate, simple and robust classification. The variables that do not

288   contain more information than the random variables will be regarded as uninformative

289   variables and eliminated using UVE. To minimize the uninformative variables the

290   random variables would be optimized with several random variables in different

291   orders of magnitude, such as $10^{-8}$, $10^{-9}$, $10^{-10}$, $10^{-11}$, $10^{-12}$. After UVE in conjunction

292   with LS-SVM, the optimal order of magnitude in random variables was $10^{-11}$. As

293   indicated in Fig.4 (a), the stability range of the added random variables with order of

294   magnitude of $10^{-11}$ (the last 1557 wavenumbers) is from 3 to -3. From the observation

295   of first 1557 variables, 700 variables whose stability lies out of the two dot lines will

296   be remained for LS-SVM models and else variables whose stability lies within the dot

297   lines will be eliminated. The remnant 700 wavenumbers were mainly located between

298   4000 and 7000cm$^{-1}$. Most of the variables between 7000 and 10000 were eliminated.

299   As a result, recognition rate of UVE-LS-SVM model achieved 96.428% shown in

300   Fig.4 (b). Only one *L. breris* was mistaken for *L. casei*. The UVE-LS-SVM model

301   based on 700 wavenumbers selected by UVE is more stable and accurate than model
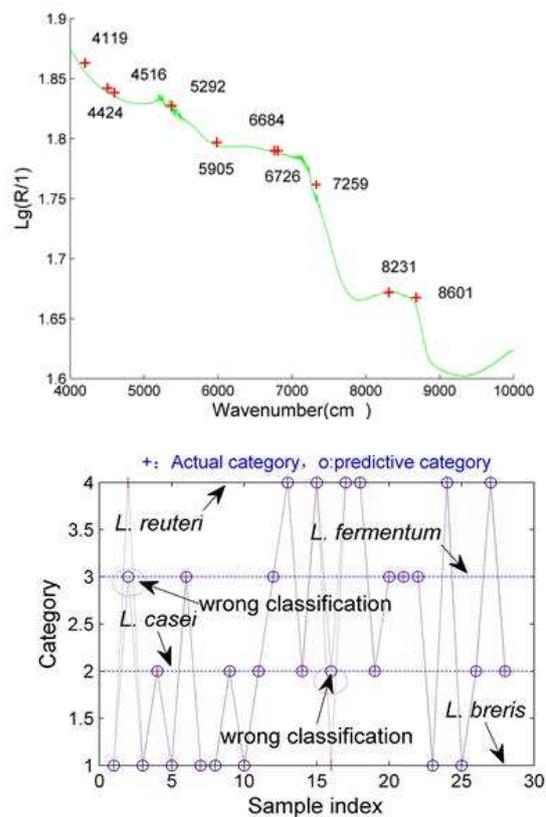
302   based on full spectrum

303       To further eliminate most of insignificant variables, GA was utilized to select the

304   characteristic wavenumbers highly correlated with the features of *Lactobacillus*. 10

305   characteristic wavenumbers (4119cm$^{-1}$, 4427cm$^{-1}$, 4317cm$^{-1}$, 5292cm$^{-1}$, 5905cm$^{-1}$,

306   6184cm$^{-1}$, 6527cm$^{-1}$, 6969cm$^{-1}$, 8531cm$^{-1}$, 8601cm$^{-1}$) were selected by GA based on

307   remaining 700 wavelengths. The 10 wavenumbers were labeled on the spectrum

308   shown in Fig.5. With the 10 characteristic wavelengths, UVE-GA-LS-SVM model

309   achieved 92.857% for the recognition rate of prediction and 100% for the recognition

310  rate of calibration. The recognition rates of UVE- GA-LS-SVM model were between

311  that of the LS-SVM model in full-spectrum and that of UVE-LS-SVM model. But

312  UVE- GA-LS-SVM with the least wavenumbers (10 wavenumbers) is simpler and

313  more robust than UVE-LS-SVM and LS-SVM models with 700 and 1557

314  wavenumbers, respectively. As shown in Fig.5 (b), the classification result of

315  UVE-GA-LS-SVM model demonstrates that one *L. breris* was mistaken for *L. casei*

316  and one *L. reuteri* was mistakenly predicted as *L. fermentum*.



317

318  Fig.4. (a): the stability distribution of each variable for identification of strains by UVE. The two dot lines in (a)

319  indicate the lower and upper threshold; (b): model performance for prediction.

320

321

322

Fig.5. (a) the ten red plus signs indicative the position of 10 characteristic wavenumbers in spectra by UVE and

GA. (b) model performance for prediction based on UVE-GA-LS-SVM.

Tab.1. overall recognition rates (%) for identification of four *Lactobacillus* based on full spectrum

LS-SVM using different pre-processing methods

| Chemometrics methods | LS-SVM | |
|---|---|---|
| | calibration | prediction |
| None | 85.714 | 82.857 |
| SNV | 98.214 | 88.571 |
| MSC | 98.214 | 91.429 |
| 1D | 89.286 | 82.857 |
| 2D | 89.286 | 82.857 |
| SNV-1D | 92.857 | 85.714 |
| SNV-2D | 92.857 | 85.714 |
| MSC-1D | 96.429 | 85.714 |
| MSC-2D | 98.214 | 85.714 |
| MSC-UVE | 98.214 | 96.428 |
| MSC-UVE-GA | 98.214 | 92.857 |

327

328

329

330    3.4 Scientific explanation of established identification models
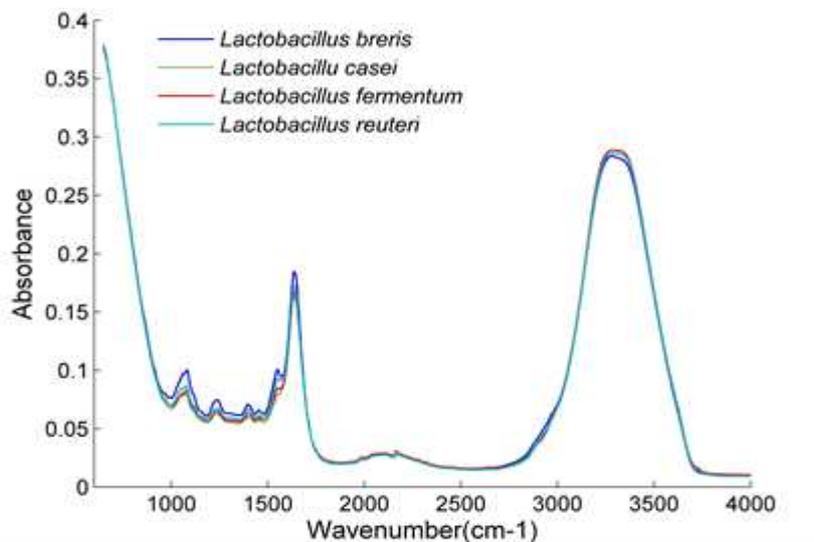
331        In this study, MIR spectroscopy was employed to investigate the components

332    and structure of bacterial cells and the absorption bands of each *Lactobacillus* was

333    used to validate the difference of components and structure among *Lactobacillus*.

334    MIR spectra of 4 bacterial cells were displayed in Fig.6a. Absorption bands of four

335    *Lactobacillus* were located at the region of 1040-1090cm$^{-1}$, 1230-1240

336    cm$^{-1}$,1390-1400, 1440-1450, 1530-1560 cm$^{-1}$, 1630-1640 cm$^{-1}$, 2850-1990 cm$^{-1}$,

337    3250-3260 cm$^{-1}$ owing to the inconsistent peaks of each spectrum. The characteristic

338    absorption in frequency at 3250-3260 cm$^{-1}$ in the frequency region of 3150

339    cm$^{-1}$-3500cm$^{-1}$ is attributed to the O-H bonded stretching vibration. 2850-1990 cm$^{-1}$

340    are in the range of 3000-2800cm-1 influenced by C-H stretching vibrations in fatty

341    acid and some amino acid. The C=O stretching vibration of amides linked to proteins

342    make contribution to the absorption peak at 1630-1640 cm$^{-1}$ (amide I) and the

343    absorption peak at 1530-1560cm$^{-1}$ is likely due to the N-H deformation of amides

344    linked to proteins (amide II). The peaks between 1500~1200 cm$^{-1}$, namely 1230-1240

345    cm$^{-1}$, 1390-1400 and 1440-1450 cm$^{-1}$ were possibly influenced by $CH_2$ and $CH_3$

346    bending modes of proteins, fatty acids and phosphate-carrying compounding. The

347    1040-1090cm$^{-1}$ between 1200 cm$^{-1}$ and 900 cm$^{-1}$ were due to the symmetric stretching

348    vibration of $PO_2^-$ groups found in nucleic acids (*44, 45*). Bacterial cells maybe all

349    contain proteins, peptides, fatty acid, and polysaccharides et al based on these

350    absorption peaks. However, intensity of peaks and the highest point of peaks among

351    the four spectra was different. That means the groups, components and the content of

352    them are different. Therefore, *Lactobacillus* could be classified based on the

353    difference of groups, components and the content of them.

354        For acquisition of MIR spectroscopy, bacterial samples, potassium bromide (KBr)

355    pellets were prepared. KBr pellets were the most common sample for acquisition of

356    MIR spectra and described by D.J.M. Mouwen et.al(*24*). About 2 mg colonies and

357    200 mg KBr powder were homogenized in an agate-stone mortar. The mixture was

358    made into a coin shaped pellet. Finally, six KBr pellets for each *Lactobacillus* were

359    obtained and used for acquisition of MIR spectra. The MIR spectra were measured

360 with Nicolet 380 FT-IR spectrometer (Thermo Electron Corporation, USA). MIR

361 spectra were acquired in the spectral range of 500 to 4000 cm$^{-1}$ at resolution of 2 cm$^{-1}$.

362 The colonies and coin shaped pellets were positioned and directly contacted with an

363 infrared attenuated total reflection diamond. Six spectra were acquired for each

364 sample and each spectrum were composed of an average of 36 separate scans. Finally,

365 six spectra were averaged. Although FT-IR has been well known as a technique for

366 identification of bacteria Bacterial samples preparation for MIR spectroscopy requires

367 a large number of bacterial cells, centrifugation, lyophilization and pressing. On

368 account of the complicated and time-consuming sample preparation, the application of

369 NIR in rapid identification of strains were restricted.

370 In this study, the NIR spectra were acquired from colonies on agar plates not the

371 KBr pellets. This sample preparation leave out most of complicated procedures

372 (centrifugation, lyophilization and pressing) except bacterial culture. The selected ten

373 characteristic wavenumbers for identification of strains were also influenced by

374 functional groups in fatty acid, peptides and proteins et al. 5905 cm$^{-1}$ is the first

375 overtone of Carbon-hydrogen (C-H); 8531 and 8601 are the weaker second overtone

376 of C-H; 6184 and 6527cm$^{-1}$ are the first overtone of Nitrogen–hydrogen (N-H); 6969

377 cm$^{-1}$ is the first overtone of Oxygen–hydrogen (O-H); 5292cm$^{-1}$ is the second

378 miscellaneous overtone band of carbonyl group in peptides; 4119, 4427 and 4317 cm$^{-1}$

379 are the C-H combination bands. The bands of C-H may be linked with fatty acid and

380 some amino acid. The bands of O-H may be linked with water. The bands of N-H

381 may be linked with proteins and peptides. The bands of miscellaneous overtone

382 maybe linked to peptides. Relationship between the ten characteristic wavenumbers in

383 NIR spectra and the absorption bands in MIR spectra was also investigated. The

384 frequency of 5905 cm$^{-1}$ is included in frequencies of two times of the region of

385 2850-2990 cm$^{-1}$. The frequencies of 8531 and 8601 cm$^{-1}$are included in frequencies of

386 three times of the region of 1210-1260cm$^{-1}$. The frequencies of 6184 and 6527 cm$^{-1}$

387 are included in frequencies of three times of the region of 1530-1560 and 1630-1640

388 cm$^{-1}$respectively. The frequency of 5292 cm$^{-1}$ is included in frequencies of three times

389 of the region of 1630-1640 cm$^{-1}$. The frequencies of 4119, 4427 and 4317 cm$^{-1}$ are the

390      combination of 1230-1240 cm$^{-1}$, 1390-1400 cm$^{-1}$, 1440-1450 and 2850-2990 cm$^{-1}$. So

391      we consider that of 5905, 8531 and 8601 cm$^{-1}$ are influenced by functional groups of

392      membrane fatty acid and by some amino acid on account of C-H stretching vibrations.

393      6184 and 6527 cm$^{-1}$ are affected by amide I and amide II groups belong to proteins

394      and peptides due to N-H stretching. 5292cm$^{-1}$ is the second miscellaneous overtone

395      band of carbonyl group in peptides owing to C=O stretching. 4119 cm$^{-1}$, 4427 cm$^{-1}$

396      and 4317 cm$^{-1}$ are the C-H combination bands influenced by vibration of C-H. The

397      NIR spectra at ten wavenumbers, influenced by the components and structure in

398      bacterial cells, are unique for each *Lactobacillus* owing to the different components

399      and structure in diverse *Lactobacillus*. Therefore, single strain could be identified

400      used the unique NIR spectra at ten wavenumbers.



401

402          Fig.6. The average IR spectra of bacterial cells (a) and bacterial colonies (b)

403      Table 2 Tentative assignment of characteristic wavenumbers found in NIR spectra

| Frequency (cm$^{-1}$) | Related frequency in MIR | Assignment |
|---|---|---|
| 5905, 8531, 8601 | 2850-2990, 1210-1260 | CH, CH$_2$, CH$_3$ |
| 6184, 6527 | 1530-1560 | Amide II band, C=O |
| 5292 | 1630-1640 | Amide I, N-H |
| 4119, 4427, 4317 | 1230-1240, 1390-1400 | Combination of C-H |

404      **4. Conclusion**

405      This paper presented a rapid method to discriminate four *Lactobacillus* strains

406    based on NIR spectroscopy technique aided by chemometric methods. The NIR

407    spectra were acquired from bacterial colonies on MRS agar medium. Rapid

408    identification of *Lactobacillus* were developed with the pre-treatment (MSC),

409    variables selection (UVE and GA), and supervised discriminant analysis (LS-SVM)

410    performed. MSC based on full wavenumber and LS-SVM showed the best

411    performance in all pre-processing methods with the recognition rate of 91.429%.

412    Utilization of UVE and GA, resulting in 10 wavenumber that have high correlation

413    with the features of *Lactobacillus*, could simplify the identification models and

414    improved the performance. The recognition rates of UVE- GA-LS-SVM model

415    (92.857%) were between that of the LS-SVM model in full-spectrum (91.429%) and

416    that of UVE-LS-SVM model (96.428%). The UVE-GA-LS-SVM model with 10

417    wavenumbers is simpler and more robust than the full-spectrum LS-SVM model with

418    1557 wavenumbers and UVE-LS-SVM model with 700 wavenumbers. By comparing

419    with absorption bands in MIR spectra, the ten characteristic wavenumbers are

420    influenced by functional groups in components and structure. Each bacteria has a

421    unique NIR spectra due to the stretching and bending vibrations of molecular bends or

422    functional groups presented in cellular components (proteins, nucleic, lipids, etc.).

423    Therefore, single strain could be identified for the unique NIR spectrum. This

424    methodology may become a powerful tool for identification of strains due to timely

425    spectral collection and high sensitivity for identification strains.

438

439

440

441

## Reference

442

443     1.   Patel, S. J., A comprehensive review on Probiotics. *Int. J. Pure App. Biosci* **2015**, *3*,
444     286-290.

445   2.   Zhai, Q.; Yin, R.; Yu, L.; Wang, G.; Tian, F.; Yu, R.; Zhao, J.; Liu, X.; Chen, Y. Q.; Zhang, H.,
446   Screening of lactic acid bacteria with potential protective effects against cadmium toxicity.
447     *Food Control* **2015**, *54*, 23-30.

448     3.   Hammes, W. P.; Vogel, R. F., The genus lactobacillus. In *The genera of lactic acid*
449     *bacteria*, Springer: 1995; pp 19-54.

450     4.   Ng, E. W.; Yeung, M.; Tong, P. S., Effects of yogurt starter cultures on the survival of
451   Lactobacillus acidophilus. *International Journal of Food Microbiology* **2011**, *145*, 169-175.

452   5.   Cogan, T. M.; Barbosa, M.; Beuvier, E.; BIANCHI-SALVADORI, B.; COCCONCELLI, P. S.;
453   FERNANDES, I.; GOMEZ, J.; GOMEZ, R.; KALANTZOPOULOS, G.; LEDDA, A., Characterization of
454     the lactic acid bacteria in artisanal dairy products. *Journal of Dairy Research* **1997**, *64*,
455     409-421.

456   6.   Wang, L.; Zhou, H.; He, R.; Xu, W.; Mai, K.; He, G., Effects of soybean meal fermentation
457   by Lactobacillus plantarum P8 on growth, immune responses, and intestinal morphology in
458     juvenile turbot (Scophthalmus maximus L.). *Aquaculture* **2016**, *464*, 87-94.

459     7.   Haghshenas, B.; Nami, Y.; Abdullah, N.; Radiah, D.; Rosli, R.; Khosroushahi, A. Y.,
460   Anticancer impacts of potentially probiotic acetic acid bacteria isolated from traditional dairy
461     microbiota. *LWT-Food Science and Technology* **2015**, *60*, 690-697.

462     8.   Waites, M. J.; Morgan, N. L.; Rockey, J. S.; Higton, G., *Industrial microbiology: an*
463     *introduction*. John Wiley & Sons: 2009.

464   9.   Geissler, A. J.; Behr, J.; von Kamp, K.; Vogel, R. F., Metabolic strategies of beer spoilage
465     lactic acid bacteria in beer. *International journal of food microbiology* **2016**, *216*, 60-68.

466     10.   García-Ruiz, A.; Crespo, J.; López-de-Luzuriaga, J.; Olmos, M.; Monge, M.;
467     Rodríguez-Alfaro, M.; Martín-Alvarez, P.; Bartolome, B.; Moreno-Arribas, M., Novel
468     biocompatible silver nanoparticles for controlling the growth of lactic acid bacteria and
469     acetic acid bacteria in wines. *Food Control* **2015**, *50*, 613-619.

470   11.   de Almeida Júnior, W. L. G.; da Silva Ferrari, Í.; de Souza, J. V.; da Silva, C. D. A.; da Costa,
471   M. M.; Dias, F. S., Characterization and evaluation of lactic acid bacteria isolated from goat
472     milk. *Food Control* **2015**, *53*, 96-103.

473   12.   Nair, P. S.; Surendran, P. K., Biochemical characterization of lactic acid bacteria isolated
474     from fish and prawn. **2005**.

13. Zhao, J.; Fleet, G., The effect of lactic acid bacteria on cocoa bean fermentation. *International journal of food microbiology* **2015,** *205*, 54-67.

14. Singh, S.; Goswami, P.; Singh, R.; Heller, K. J., Application of molecular identification tools for Lactobacillus, with a focus on discrimination between closely related species: a review. *LWT-Food Science and Technology* **2009,** *42*, 448-457.

15. Sohier, D.; Coulon, J.; Lonvaud-Funel, A., Molecular identification of Lactobacillus hilgardii and genetic relatedness with Lactobacillus brevis. *International Journal of Systematic and Evolutionary Microbiology* **1999,** *49*, 1075-1081.

16. Coeuret, V.; Dubernet, S.; Bernardeau, M.; Gueguen, M.; Vernoux, J. P., Isolation, characterisation and identification of lactobacilli focusing mainly on cheeses and other dairy products. *Le Lait* **2003,** *83*, 269-306.

17. Montville, T. J.; Matthews, K. R., *Food microbiology: an introduction*. ASM Press: 2005.

18. Moreira, J. L. S.; Mota, R. M.; Horta, M. F.; Teixeira, S. M.; Neumann, E.; Nicoli, J. R.; Nunes, Á. C., Identification to the species level of Lactobacillus isolated in probiotic prospecting studies of human, animal or food origin by 16S-23S rRNA restriction profiling. *BMC microbiology* **2005,** *5*, 1.

19. Adeyemo, S.; Onilude, A., Molecular identification of Lactobacillus plantarum isolated from fermenting cereals. *International Journal of Biotechnology and Molecular Biology Research* **2014,** *5*, 59-67.

20. Anukam K C, O. E. O., Ahonkhai I, 16S rRNA gene sequence and phylogenetic tree of Lactobacillus species from the vagina of healthy Nigerian women. *African Journal of Biotechnology* **2005,** *4*, 1222-1227.

21. Tenover F C, A. R. D., Goering R V, Interpreting chromosomal DNA restriction patterns produced by pulsed-field gel electrophoresis: criteria for bacterial strain typing. *Journal of clinical microbiology* **1995,** *33*, 2233.

22. Bidet, P., Barbut, F., Lalande, V., Burghoffer, B., & Petit, J. C., Development of a new PCR-ribotyping method for Clostridium difficile based on ribosomal RNA gene sequencing. *FEMS microbiology letters* **1999,** *175*, 261-266.

23. Marty, E.; Buchs, J.; Eugster-Meier, E.; Lacroix, C.; Meile, L., Identification of staphylococci and dominant lactic acid bacteria in spontaneously fermented Swiss meat products using PCR–RFLP. *Food microbiology* **2012,** *29*, 157-166.

24. Mouwen, D.; Hörman, A.; Korkeala, H.; Alvarez-Ordóñez, A.; Prieto, M., Applying Fourier-transform infrared spectroscopy and chemometrics to the characterization and identification of lactic acid bacteria. *Vibrational Spectroscopy* **2011,** *56*, 193-201.

25. Grasso, E. M.; Yousef, A. E.; de Lamo Castellvi, S.; Rodriguez-Saona, L. E., Rapid detection and differentiation of Alicyclobacillus species in fruit juice using hydrophobic grid membranes and attenuated total reflectance infrared microspectroscopy. *Journal of agricultural and food chemistry* **2009,** *57*, 10670-10674.

26. Lecellier, A.; Gaydou, V.; Mounier, J.; Hermet, A.; Castrec, L.; Barbier, G.; Ablain, W.; Manfait, M.; Toubas, D.; Sockalingum, G., Implementation of an FTIR spectral library of 486 filamentous fungi strains for rapid identification of molds. *Food microbiology* **2015,** *45*, 126-134.

517    27.  Driver, T.; Bajhaiya, A. K.; Allwood, J. W.; Goodacre, R.; Pittman, J. K.; Dean, A. P.,

518    Metabolic responses of eukaryotic microalgae to environmental stress limit the ability of

519    FT-IR spectroscopy for species identification. *Algal research* **2015**, *11*, 148-155.

520  28.  Alvarez-Ordóñez, A.; Mouwen, D. J. M.; López, M.; Prieto, M., Fourier transform infrared

521    spectroscopy as a tool to characterize molecular composition and stress response in

522    foodborne pathogenic bacteria. *Journal of Microbiological Methods* **2011,** *84*, 369-378.

523  29.  Bosch, A.; Golowczyc, M. A.; Abraham, A. G.; Garrote, G. L.; De Antoni, G. L.; Yantorno,

524    O., Rapid discrimination of lactobacilli isolated from kefir grains by FT-IR spectroscopy.

525    *International Journal of Food Microbiology* **2006,** *111*, 280-287.

526    30.  Feng, Y.-Z.; Downey, G.; Sun, D.-W.; Walsh, D.; Xu, J.-L., Towards improvement in

527  classification of Escherichia coli, Listeria innocua and their strains in isolated systems based

528    on chemometric analysis of visible and near-infrared spectroscopic data. *Journal of Food

529    Engineering* **2015,** *149*, 87-96.

530    31.  Oust, A.; Møretrø, T.; Kirschner, C.; Narvhus, J. A.; Kohler, A., FT-IR spectroscopy for

531    identification of closely related lactobacilli. *Journal of Microbiological Methods* **2004,** *59*,

532    149-162.

533  32.  Oelofse, A.; Malherbe, S.; Pretorius, I. S.; Du Toit, M., Preliminary evaluation of infrared

534  spectroscopy for the differentiation of Brettanomyces bruxellensis strains isolated from red

535    wines. *International Journal of Food Microbiology* **2010,** *143*, 136-142.

536    33.  Marques, A. S.; Castro, J. N.; Costa, F. J.; Neto, R. M.; Lima, K. M., Near-infrared

537    spectroscopy and variable selection techniques to discriminate Pseudomonas aeruginosa

538    strains in clinical samples. *Microchemical Journal* **2016,** *124*, 306-310.

539  34.  Rodriguez-Saona, L. E.; Khambaty, F. M.; Fry, F. S., &; Calvey, E. M., Rapid detection and

540  identification of bacterial strains by Fourier transform near-infrared spectroscopy. *Journal of

541    agricultural and food chemistry* **2001,** *49*, 574-579.

542    35.  Rinnan, Å.; Berg, F. v. d.; Engelsen, S. B., Review of the most common pre-processing

543    techniques for near-infrared spectra. *TrAC Trends in Analytical Chemistry* **2009,** *28*,

544    1201-1222.

545    36.  Guo, Y.; Ni, Y.; Kokot, S., Evaluation of chemical components and properties of the

546  jujube fruit using near infrared spectroscopy and chemometrics. *Spectrochimica Acta Part A:

547    Molecular and Biomolecular Spectroscopy* **2016,** *153*, 79-86.

548    37.  Berrueta, L. A.; Alonso-Salces, R. M.; Héberger, K., Supervised pattern recognition in

549    food analysis. *Journal of Chromatography A* **2007,** *1158*, 196-214.

550    38.  Mordehai, J.; Ramesh, J.; Huleihel, M.; Cohen, Z.; Kleiner, O.; Talyshinsky, M.;

551  Erukhimovitch, V.; Cahana, A.; Salman, A.; Sahu, R. K., Studies on acute human infections

552    using FTIR microspectroscopy and cluster analysis. *Biopolymers* **2004,** *73*, 494-502.

553  39.  Xie, C.; Mace, J.; Dinno, M.; Li, Y.; Tang, W.; Newton, R.; Gemperline, P., Identification of

554  single bacterial cells in aqueous solution using confocal laser tweezers Raman spectroscopy.

555    *Analytical chemistry* **2005,** *77*, 4390-4397.

556    40.  Centner, V.; Massart, D.-L.; de Noord, O. E.; de Jong, S.; Vandeginste, B. M.; Sterna, C.,

557  Elimination of uninformative variables for multivariate calibration. *Analytical chemistry* **1996,**

558    *68*, 3851-3858.

559    41.  de Sousa Marques, A.; de Melo, M. C. N.; Cidral, T. A.; de Lima, K. M. G., Feature

560    selection strategies for identification of Staphylococcus aureus recovered in blood cultures

561    using FT-IR spectroscopy successive projections algorithm for variable selection: a case study. *Journal of microbiological methods* **2014,** *98*, 26-30.

563    42.  Pavlova, S. I.; Kilic, A.; Kilic, S.; So, J. S.; Nader‑Macias, M.; Simoes, J.; Tao, L., Genetic diversity of vaginal lactobacilli from women in different countries based on 16S rRNA gene sequences. *Journal of Applied Microbiology* **2002,** *92*, 451-459.

566    43.  Song, Y.-L.; Kato, N.; Liu, C.-X.; Matsumiya, Y.; Kato, H.; Watanabe, K., Rapid identification of 11 human intestinal Lactobacillus species by multiplex PCR assays using group-and species-specific primers derived from the 16S–23S rRNA intergenic spacer region and its flanking 23S rRNA. *FEMS Microbiology Letters* **2000,** *187*, 167-173.

570    44.  Alvarez-Ordonez, A.; Mouwen, D.; Lopez, M.; Prieto, M., Fourier transform infrared spectroscopy as a tool to characterize molecular composition and stress response in foodborne pathogenic bacteria. *Journal of microbiological methods* **2011,** *84*, 369-378.

573    45.  Man, Y. C.; Mirghani, M. E. S., Rapid method for determining moisture content in crude palm oil by Fourier transform infrared spectroscopy. *Journal of the American Oil Chemists' Society* **2000,** *77*, 631-637.