

Analytical Methods

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

1 Identification of heavy metal-contaminated *Tegillarca* 2 *granosa* using infrared spectroscopy

3 Xiaojing Chen^a, KeLiu^a, JingboCai^c, DehuaZhu^b, Huiling Chen^{a*}

4 ^a College of Physics and Electronic Engineering Information, Wenzhou University

5 ^b College of Mechanical& Electrical Engineering, Wenzhou University

6 ^c Zhejiang Key Laboratory of Exploitation and Preservation of Coastal Bio-resource, Zhejiang

7 Mariculture Research Institute

8 * Corresponding author. Tel.: +86 577 86689027; fax +86 577 86689027

9 E-mail: Huiling@wzu.edu.cn (Huiling Chen*)

10 **Abstract:** This study explored the feasibility of using infrared spectroscopy for the rapid detection
11 of heavy metal contamination in *Tegillarca granosa*. Generally, there is no specific characteristic
12 peak of heavy metals in the infrared range. Nevertheless, these are some changes in the structure
13 and concentration of relevant biological molecules induced by heavy metal contamination produce
14 remarkably weak spectral information on heavy metals. In this study, we selected characteristic
15 infrared spectral variables to obtain heavy metal information using the Competitive Adaptive
16 Reweighted Sampling method, Successive Projection Algorithm and Genetic Algorithm. The
17 selected variables served as inputs for classification algorithm to construct two classification
18 models. One model was designed to classify *Tegillarca granosa* samples that were
19 uncontaminated (healthy) and contaminated by a certain heavy metal (Cu, Cd, Pb, or Zn)
20 (Design I). The other model was designed to classify all sample varieties, including the samples
21 that were uncontaminated and contaminated by the four heavy metals (Design II). The two models

1
2
3
4 were validated using 10-fold cross validation. The prediction accuracy by combination of
5
6 Competitive Adaptive Reweighted Sampling method and Support Vector Machine algorithm
7
8 reached 95% for Design I and 92% for Design II. The results of this study indicated the potential
9
10 of infrared spectroscopy in evaluating heavy metal contamination in *Tegillarca granosa*.
11
12

13 **Keywords:** Heavy metal; Infrared spectroscopy; Aquatic product; *Tegillarca granosa*; CARS;
14
15 Support vector machine;
16
17

18 **Introduction**

19
20 *Tegillarca granosa* is a nutritious food source because of its low cholesterol content and high
21
22 protein, ferrum, calcium, carbohydrate, riboflavin, and various trace element contents [1-2].
23
24 However, *Tegillarca granosa* principally thrives in tidelands close to sewages, which are highly
25
26 exposed to heavy metal contamination. Given its non-selective filter-feeding behavior and low
27
28 mobility, *Tegillarca granosa* can accumulate heavy metals at concentrations ten to thousand times
29
30 higher than other aquatic species. The accumulated heavy metals remain in *Tegillarca granosa* for
31
32 a long time, making this seafood hazardous to human health. Soluble heavy metal ions are
33
34 absorbed through the gills of *Tegillarca granosa* and then distributed throughout the body through
35
36 blood circulation. These ions can be accumulated in specific body parts or on surface cells; they
37
38 can also be absorbed through the digestive tract during feeding. *Tegillarca granosa* absorbs
39
40 soluble and particulate heavy metals. Soluble heavy metals are principally absorbed by the body
41
42 surface, whereas particulate heavy metals are absorbed through the ingestion and digestion of
43
44 heavy metal-contaminated food [3]. Upon entering the food chain, these heavy metals accumulate
45
46 to high concentrations through bio-concentration and bio-magnification. Some metals can be
47
48 easily discharged, whereas others may accumulate and affect *Tegillarca granosa* tissues. When
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 44 ingested, *Tegillarca granosa* containing toxic levels of heavy metals can cause serious harm to
5
6 45 human health. Therefore, it is important and necessary to detect heavy metal pollution in
7
8
9 46 *Tegillarca granosais*.

10
11 47 At present, common methods of heavy metal detection include graphite furnace atomic
12
13 48 absorption spectrometry, flame atomic absorption spectrometry, atomic fluorescence spectrometry,
14
15 49 and inductively coupled plasma mass spectrometry [4-5]. However, these methods are expensive,
16
17 50 labor intensive, complex, time consuming, and require a large sample size. Therefore, a fast,
18
19 51 simple and reliable method that overcomes all these drawbacks is necessary to be developed to
20
21 52 detect heavy metal contamination in *Tegillarca granosa*. In addition, *Tegillarca granosa* can
22
23 53 reflect the level of heavy metal contamination in its surrounding environment. *Tegillarca granosa*,
24
25 54 oyster, and other shellfish are widely used as indicator organisms of heavy metal contamination
26
27 55 for the biological monitoring of marine pollution in numerous countries. Therefore, rapid
28
29 56 detection of heavy metal contamination in *Tegillarca granosais* urgently needed.

30
31 57 Infrared (IR) spectroscopy has been widely used in food safety and quality control because of
32
33 58 its rapidity, simplicity, high precision, low maintenance cost and small sample size requirement
34
35 59 [6-10]. IR spectroscopy detects the vibrational and rotational energies of molecules within the IR
36
37 60 spectrum by measuring the absorbance spectrum of hydrogen bonding. This technology also
38
39 61 gathers useful information that reflects organic molecules such as proteins, lipids and sugar in
40
41 62 biological tissues. IR spectroscopy is an important strategy for the structural analysis of organic
42
43 63 compounds [11-12]. However, heavy metals generally do not any show IR activity and barely
44
45 64 have any characteristic peak in the IR spectrum. Heavy metals only indirectly change the
46
47 65 vibrational spectrum by inhibiting antioxidant enzymes or by inducing the synthesis of
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 66 detoxification proteins in large quantities. These phenomena change the structures and
5
6 67 concentrations of relevant biological molecules [11, 13-17]. Therefore, IR spectral information on
7
8
9 68 heavy metal contamination can be indirectly obtained through the interactions between heavy
10
11 69 metal ions and enzymes. However, this information is remarkably weak. An in-depth analysis of
12
13
14 70 these weak signals is therefore crucial to detect heavy metal accumulation in *Tegillarca granosa*.

15
16 71 This study aims to select characteristic spectral variables that differentiate *Tegillarca granosa*
17
18 72 samples that were uncontaminated (healthy) from those that were contaminated by copper (Cu),
19
20
21 73 cadmium (Cd), zinc (Zn), and lead (Pb) using the Competitive Adaptive Reweighted
22
23
24 74 Sampling(CARS) method, Successive Projection Algorithm (SPA) and Genetic Algorithm (GA).
25
26 75 Discrimination models that classify healthy samples and heavy metal-contaminated samples were
27
28
29 76 constructed using the selected spectral variables. Cd and Pb are not essential and toxic, whereas
30
31 77 Zn and Cu are crucial to the human body. However, excessive Zn and Cu can also harm human
32
33
34 78 health. Hence, detecting heavy metal contamination in *Tegillarca granosa* is important to ensure
35
36 79 food safety. The specific objectives of the current work were to (1) spectral variables that
37
38
39 80 differentiate *Tegillarca granosa* samples that were uncontaminated (healthy) and contaminated by
40
41 81 a certain heavy metal were selected using CARS, SPA and GA (Design I). classification models
42
43
44 82 were then constructed using these selected spectral variables; (2) Spectral variables that
45
46 83 differentiate all sample varieties, including the samples that were uncontaminated and
47
48
49 84 contaminated by any of the four heavy metals, were selected using CARS, SPA and GA (Design
50
51 85 II). Classification model was then constructed with these selected spectral variables.

86 **Materials and methods**

87 **Sample Preparation and Chemicals**

1
2
3
4 88 *Tegillarca granosa* samples were purchased from Zhejiang Mariculture Research Institute at
5
6 89 Wenzhou, China in May 2014, and were acclimatized to laboratory conditions for approximately
7
8
9 90 10 d in plastic pools with a size of 60.0 cm×40.0 cm×30.0 cm. Analytical-grade
10
11 91 $\text{PbCH}_3\text{COO}\cdot 3\text{H}_2\text{O}$, $\text{CuSO}_4\cdot 5\text{H}_2\text{O}$, CdCl_2 , and $\text{ZnSO}_4\cdot 7\text{H}_2\text{O}$ were purchased from Chemical
12
13
14 92 Reagent Co. Ltd., Shanghai, China.

15
16 93 Seawater prepared by over 24h of sedimentation and then sand filtration was used to
17
18 94 maintain *Tegillarca granosa* in the tanks. The prepared seawater had a pH of 8.05 ± 0.1 , a
19
20
21 95 temperature of 20.8 ± 2.6 °C, a dissolved oxygen content of > 6 mg/L, and a salinity level of 21‰.
22
23
24 96 The water was changed every 24h throughout the experiment. The containers were refilled and
25
26 97 re-dosed with the metal toxicant.

27
28
29 98 The *Tegillarca granosa* samples in Groups I, II, III, and IV were exposed to high
30
31 99 concentrations of $\text{PbCH}_3\text{COO}\cdot 3\text{H}_2\text{O}$ (1.833 mg/L), $\text{CuSO}_4\cdot 5\text{H}_2\text{O}$ (5.589mg/L), CdCl_2 (1.634
32
33
34 100 mg/L), and $\text{ZnSO}_4\cdot 7\text{H}_2\text{O}$ (4.424 mg/L) in water, respectively. Group V (control) was reared in
35
36 101 seawater without adding heavy metals. The *Tegillarca granosa* samples from all groups were reared
37
38 102 for 10 d to allow heavy metal accumulation. After the rearing period, the *Tegillarca granosa*
39
40
41 103 samples were sacrificed and then stored in a refrigerator at -4 °C for 15 min. The samples were
42
43
44 104 freeze-dried, ground into powder, and then used for spectral analyses.

45
46 105 A total of 150 samples (30 samples for each variety) were prepared for further treatment, in
47
48 106 which 30 samples from the healthy group and 30 samples from the contaminated group were used
49
50
51 107 to establish the models for Designs I and II. Therefore, 60 samples in four sample sets were used
52
53
54 108 for Design I, and these four sample sets were used for the intoxication analysis of Zn, Cu, Cd, and
55
56 109 Pb. Meanwhile, 150 samples (30 samples×5 groups) were used for Design II. A column vector (Y_C)

1
2
3
4 110 containing the integer numbers of intoxication status of samples from the training set was
5
6 111 concatenated to matrix X_C . Model training was based on X_C and Y_C . For the two-class problem
7
8 112 (Design I), the integer numbers were set to 0 and 1, which represent the contaminated and healthy
9
10 113 samples, respectively. For the five-class problem (Design II), the integer numbers were set to 0
11
12 114 and 1 to 4, which represent the healthy samples and the samples contaminated by Cd, Cu, Pb, and
13
14
15 115 Zn, respectively.

116 **Spectral Collection and Reference Methods for Heavy Metal Content**

117 IR (4000–400 cm^{-1}) spectra were obtained using a Tensor 27 spectrophotometer (Bruker, Inc.,
118 Germany) equipped with a Golden Gate Diamond ATR sampling accessory. The collection of all
119 samples was completed in an airtight collection box. The instrument need warm up about 30
120 minutes before measuring the spectrum. All samples were scanned 15 times, in which a scan time
121 is about 5 seconds, and the results were averaged using OMNIC software (Version 5.2, Bruker,
122 Inc.). The SNR of spectrophotometer is superior to 55000:1 (peak-to-peak value), resolution ration
123 is 4cm^{-1} , KBr beam splitter, DTGS detector and 100-micrometer diaphragm were adopted.

124 The reference value of the heavy metal concentration was measured using a NexION 300X
125 ICP-MS (NexION 300X, Perkin Elmer, Inc., U.S.). The gas flow rate, auxiliary gas flow rate, RF
126 power, and peristaltic pump of the atomized sample were set to 0.90L/min, 14L/min, 1100W, and
127 20rpm, respectively. The detailed test procedure was as follows. First, exactly 0.3g of sample
128 (deviation < 0.001 g) was prepared, digested with 6mL of highly pure nitric acid, and then filled to
129 50mL with distilled water. The sample was subjected to microwave digestion, heated to 120 °C
130 within 5min, and then maintained at 120 °C for 10min. The sample was heated from 120 °C to
131 180 °C for 10min and then maintained at 180 °C for 20min. Afterward, the sample was analyzed

1
2
3
4 132 using the instrument tested. The statistical values of the heavy metal contents of the samples are
5
6 133 shown in Table 1. The table shows that the four heavy metals were mildly enriched in the healthy
7
8
9 134 samples but were highly enriched in the contaminated samples.

11 135 **Variable Selection method**

12
13
14 136 CARS method

15
16 137 Competitive Adaptive Reweighted Sampling (CARS) is a newly proposed spectral variable
17
18 138 selection method whose algorithm is based on the “survival of the fittest” principle of Darwin’s
19
20
21 139 theory of evolution [18-21]. The key concept is to consider each spectral variable as an individual
22
23
24 140 and to remove unfit individuals. Spectral variables with large absolute coefficients in partial least
25
26 141 squares (PLS) regression were selected and those with small values were excluded using CARS.
27
28
29 142 Spectral variable subsets were obtained, and each subset was used to construct a model for CV.
30
31 143 The model with the lowest root mean square error of CV (RMSECV) was selected as the optimal
32
33
34 144 subset of wavelengths. The model with the lowest RMSECV effectively selected the optimized
35
36 145 subset of spectral variables relevant to properties of interest. The algorithm was as follows:

37
38
39 146 (1) The population was subjected to N rounds of Monte Carlo (MC) sampling. In each
40
41 147 sampling, a PLS regression model was constructed with a random partition of samples
42
43
44 148 selected for validation. The absolute value of the regression coefficient $|b_i|$ was calculated,
45
46 149 where i represents the i^{th} model.

47
48
49 150 (2) Spectral variables with relatively small $|b_i|$ were forcibly removed using the exponential
50
51 151 decreasing function:

$$122 \quad r_i = ae^{-ki} \quad (1)$$

53
54
55
56 152 where the constants a and k were calculated as follows:
57
58
59
60

1
2
3
4 154 $a = (p/2)^{1/(N-1)}$ (2)
5

6 155 $k = \ln(p/2)/(N-1)$ (3)
7

8
9 156 In the first round of MC sampling, all p variables were used for modeling; thus, $r_i = 1$, In
10
11 157 the N^{th} MC sampling, only two spectral variables were used; thus, $r_N = p/2$.

12
13
14 158 (3) N rounds of CARS selection were performed. In each round, the spectral variables with
15
16 159 large absolute values of PLS regression coefficients were selected and used to construct a
17
18 160 PLS regression model. Then, the spectral variables with the lowest RMSECV values were
19
20 161 selected.

21
22
23
24 162 SPA method

25
26 163 Successive projections algorithm (SPA) is a forward variable selection method in the area of
27
28 164 spectral analysis. It uses simple operations in a vector space to select a subset of variables whose
29
30 165 information content is minimally redundant for multivariate linear regression (MLR) [22-23]. In
31
32 166 SPA calculation, candidate subsets of variables are constructed according to a sequence of
33
34 167 projection operations involving the columns of the instrumental response matrix. These candidate
35
36 168 subsets are evaluated according to the prediction performance of the resulting MLR model. In SPA,
37
38 169 such prediction performance is assessed by an independent validation set or cross-validation of
39
40 170 training set. With the advantages of simple, stability and good predictive performance, SPA has
41
42 171 been successfully employed for variable selection in UV-VIS, ICP-AES and NIR spectrometry, as
43
44 172 well as for coefficient selection in wavelet regression model.[24-27].
45
46
47
48

49
50
51 173 GA method

52
53
54 174 The ultimate goal of genetic algorithm (GAs) is the optimization of a given response function.

55
56 175 GAs are inspired by evolution theory: in a living environment, the “best” individuals have a
57
58
59
60

1
2
3
4 176 greater chance to survive and a greater probability to spread their genomes by reproduction,
5
6 177 namely the 'struggle for life'. In wavelength variables selection, GAs have five basic steps,
7
8
9 178 variable coding, population initiation, response evaluation, reproductions, and mutations. Steps
10
11 179 3-5 alternate until a termination criterion is reached. This criterion can be based on a lack of
12
13
14 180 improvement in the response or simply on a maximum number of generations or on the total time
15
16 181 allowed for elaboration. [28-29].
17

182 **Classification models**

183 PLS-DA method

184 Partial least-squares discriminant analysis (PLS-DA) is the classification version of partial
185 least-squares regression (PLSR). Differing from Partial least-squares regression (PLSR) that is
186 PLSR of a set Y of binary variables describing the categories of a categorical variable on a set X
187 of predictor variables. It is a compromise between the usual discriminant analysis and a
188 discriminant analysis on the significant principal components of the predictor variables PLS-DA
189 encodes the dependent variable with dummy variables describing the classes for the optimum
190 separation of classes. After encoding, PLS-DA is implemented in the usual way of PLSR [30].

191 SIMCA method

192 Soft independent modelling of class analogy (SIMCA) is an established method for
193 multivariate classification. In SIMCA classification, the residuals of several disjoint PCA models
194 are utilized to assign an observation to one or several of the available classes. During the training
195 of each class-specific PCA model, a distribution of the residuals for each class is generated. Given
196 this class-specific residual distribution, any given observation can subsequently be assigned a
197 probability of equal variance compared to the model residuals according to a *F*-test. The

1
2
3
4 198 probability assignment is then ultimately used to accept or reject the observation to or from each
5
6 199 class, which is essentially a tool for detecting model outliers [31-32].
7

8
9 200 **LDA method**

10
11 201 Linear discriminant analysis (LDA) is one commonly used technique for data classification
12
13 202 and dimensionality reduction. LDA is a mathematical transformation method from
14
15 203 multidimensional space to one-dimensional space by maximizing the Fisher criterion, namely the
16
17 204 ratio of between-class scatter dispersion and within-class scatter dispersion, to find the best and
18
19 205 most easily categorized projection direction. LDA easily handles the case where the within-class
20
21 206 frequencies are unequal and their performance has been examined on randomly generated test data.
22
23 207 This method maximizes the ratio of between-class variance to the within-class variance in any
24
25 208 particular data set thereby guaranteeing maximal separability. The use of Linear Discriminant
26
27 209 Analysis for data classification is applied to classification problem in speech recognition. When
28
29 210 transformed to a different space, LDA doesn't change the location of the original data sets but
30
31 211 only tries to provide more class separability and draw a decision region between the given classes.
32
33 212 This method also helps to better understand the distribution of the feature data. In LDA, data sets
34
35 213 can be transformed and test vectors can be classified in the transformed space by two different
36
37 214 approaches: class-dependent transformation and class-independent transformation [33].
38
39
40
41
42
43
44

45
46 215 **SVM method**

47
48
49 216 This section briefly describes Support Vector Machine (SVM). For further details, one can
50
51 217 refer to [34-35], which provides a complete description of the SVM theory. The linear SVM finds
52
53 218 an optimal separating margin by solving the following optimization task:
54
55
56
57
58
59
60

$$\begin{aligned}
 & \text{Minimize } g(\mathbf{w}, \xi) = \frac{1}{2} \|\mathbf{w}\|^2 + \kappa \sum_{i=1}^n \xi_i \\
 & \text{st: } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0
 \end{aligned} \quad (4)$$

where κ is a penalty value and ξ_i is the positive slack variables. This primal problem can be reduced to the Lagrangian dual problem by introducing Lagrangian multipliers (α_i). The optimal solution α_i can be obtained under the Karush–Kuhn–Tucker conditions. If $\alpha_i > 0$, the corresponding data points are called SVs. Afterward, the optimal hyperplane can be constructed using the optimal parameters w and b . The linear classification function can then be given by

$$g(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \right) \quad (5)$$

The original input space was mapped into a high-dimensional feature space via a mapping function to maximize the function of the linear learning machine in non-linear cases. Through this mapping function, $\mathbf{x}^T \mathbf{x}$ in the input space is a form of $\phi(\mathbf{x}_i)^T \phi(\mathbf{x})$ in the feature space. The Gaussian kernel is generally utilized to detect the optimal parameter values of the radial basis function (RBF) kernel (i.e., C and γ). Therefore, the decision function can be expressed as follows:

$$g(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right) \quad (6)$$

In the SVM algorithm, the RBF kernel is utilized to map the input space into a high-dimensional feature space. Thus, the parameters C and γ of RBF kernel are important in designing an effective SVM model. The penalty parameter C determines the trade-off between the fitting error minimization and model complexity, while the kernel width γ defines the non-linear mapping from the input space to a high-dimensional feature space [36].

237 Model Evaluation and Algorithm Parameters

238 The sensitivity, specificity, and accuracy of the classification models for training and

1
2
3
4 239 prediction were evaluated. Sensitivity is defined as the number of positives (contaminated samples)
5
6 240 correctly classified by the model divided by the number of all positives. Specificity is defined as
7
8 241 the number of negatives (healthy samples) correctly identified by the model divided by the
9
10 242 number of all negatives. Accuracy is defined as the number of correctly distinguished samples
11
12 243 divided by the number of all samples. All models were validated through 10-fold CV. All
13
14 244 computations and chemometric analyses were operated in Matlab 2011a (The Mathworks, Inc.,
15
16
17
18
19 245 Natick, MA, USA).

22 246 **Results and Discussion**

25 247 There were large of noise interference existed in $623\text{--}400\text{ cm}^{-1}$ and $4000\text{--}3517\text{ cm}^{-1}$ bands,
26
27 248 therefore these two bands were removed, resulting in only the band of $3516\text{--}624\text{ cm}^{-1}$ was used for
28
29 249 further analysis. Figure 1 shows a typical IR spectrum of *Tegillarca granosa*. From Figure 1,
30
31 250 healthy samples and samples contaminated by different heavy metals show similar IR spectra.
32
33 251 Each variety only has one sample, the IR spectra in the Figure 1 don't represent the true
34
35 252 differences among all samples. Considering all samples produces overlapping spectral profiles. As
36
37 253 a result, evident difference among all samples can hardly be detected by naked eye. In addition,
38
39 254 the IR spectrum of *Tegillarca granosais* complicated and contains numerous absorption peaks
40
41 255 contributed by different functional groups. The wavenumber assignment mainly concentrates in
42
43 256 two distinct ranges, namely, $3700\text{--}2800$ and $1800\text{--}650\text{ cm}^{-1}$. Comprehensively, the broad band
44
45 257 centered at around 3300 cm^{-1} is attributed to the N-H stretching mode of Amide A [11, 14-15]. In
46
47 258 the region between 3000 and 2800 cm^{-1} , three weak absorption peaks are found at around 2960 ,
48
49 259 2925 and 2865 cm^{-1} . The absorption peaks at 2960 and 2865 cm^{-1} are attributed to the asymmetric
50
51
52
53
54
55
56
57
58 260 and symmetric stretching vibrations of CH_3 . These two peaks are often used to determine the lipid
59
60

1
2
3
4 261 structure and monitor proteins in biological systems. Another absorption band at 2925 cm^{-1}
5
6 262 corresponds to the anti symmetric and symmetric stretching vibrations of $-\text{CH}_2$, which is mainly
7
8
9 263 found in lipids [11, 14-15]. Several absorbance peaks are located in the band region range of
10
11 264 $1800\text{--}650\text{ cm}^{-1}$. This region is dominated by amide groups. Two strong absorption bands at
12
13
14 265 around 1650 and 1540 cm^{-1} correspond to the amide I and amide II vibrations of structural
15
16 266 proteins, respectively. Amide I absorption is principally related to the $\text{C}=\text{O}$ stretching vibration of
17
18
19 267 amides. Meanwhile, amide II absorption originates from amide N-H bending vibration (60%)
20
21 268 coupled with the C-N stretching vibration (40%) mode of the polypeptide and protein backbones
22
23
24 269 [13]. These two absorption bands are sensitive and can be used to determine the secondary
25
26 270 structure of proteins. The absorption bands at around 1400 cm^{-1} are assigned to COO^- symmetric
27
28
29 271 stretching modes, which are mainly associated with fatty acids and amino acids [13-14]. The
30
31
32 272 bands observed at 1230 cm^{-1} correspond to the PO_2 asymmetric stretching of nucleic acids in
33
34 273 phospholipids. The band observed at 1070 cm^{-1} corresponds to the symmetric stretching of PO_2^-
35
36 274 in nucleic acids, HO-C-H stretch, and carbohydrates. The band observed at 1040 cm^{-1} is assigned
37
38
39 275 to the C-O stretching vibrations in polysaccharides. In addition, the peaks at $1000\text{--}623\text{ cm}^{-1}$
40
41
42 276 correspond to a fingerprint region mostly of nucleic acids [15]. We also observed the spectra of the
43
44 277 chemical molecules and functional groups of the constituents found in *Tegillarca granosa*. This
45
46 278 spectral information is important to detect heavy metal contamination. However, the *Tegillarca*
47
48
49 279 *granosa* samples have numerous absorption peaks that overlap between different samples. This
50
51
52 280 overlapping complicates the selection of useful spectral variables for the detection of heavy metal
53
54 281 contamination. Therefore, variable feature selection and pattern recognition methods were utilized.
55
56 282 In our study, the preprocessed method of standard normal variate (SNV) was employed to achieve
57
58
59
60

1
2
3
4 283 a centering and scaling effect. The processed spectra were used for further data analysis.
5

6
7 **284 Analysis of classification Results based on different variable selection methods**
8

9 285 Because there are lots of irrelevant information and redundant information in infrared spectra,
10
11 286 the common variable selection methods of CARS, SPA and GA were employed to select the
12
13 287 characteristic spectral variables. The selected variables by CARS, SPA and GA were served as
14
15 288 input variables to construct classification model including PLS-DA, SIMCA, LDA and SVM, so
16
17 289 48 models for Design I (four heavy metals \times four calibration methods \times three variable selection
18
19 290 methods) and 12 models for Design II (four calibration methods \times three variable selection
20
21 291 methods) was obtained. For Design I, classification results of PLS-DA, LDA and SIMCA models
22
23 292 are below than 90% except SVM model, and classification results obtained by these linear models
24
25 293 (PLS-DA, LDA and SIMCA) are only close to 70% for Design II, which are unexceptionally
26
27 294 inferior to the results of SVM which classification results exceed 90% . So, in our manuscript, the
28
29 295 SVM model as main model was employed to further analyze.
30
31
32
33
34
35

36 296 As mentioned above, the parameters C and γ of RBF kernel are important in designing an
37
38 297 effective SVM model. In this study, grid search was utilized using 10-fold CV to optimize the two
39
40 298 key parameters of the RBF kernel-based SVM. The values of C and γ vary as $C = \{2^{-5}, 2^{-3}, \dots, 2^{15}\}$
41
42 299 and $\gamma = \{2^{-15}, 2^{-13}, \dots, 2^1\}$. The optimal parameter pair (C, γ) was used to construct a model for
43
44 300 training. The classification results for Design I are shown in Table 2 and those of Design II in
45
46 301 Table 3. From Tables 2 and 3, it can be found that the variable number selected by SPA is less than
47
48 302 those of CARS and GA, meanwhile the results of GA and SPA are slightly inferior to that of
49
50 303 CARS, that might be because it is easy for GA to get partially optimized, in most cases the
51
52 304 selected variables are not the optimal. As the number of variables selected by SPA must be less
53
54
55
56
57
58
59
60

1
2
3
4 305 than that of samples involved in the modeling and there were a few samples considered in this
5
6 306 study, only some useful information (variable number was less than sample number) can be
7
8
9 307 selected by SPA. As to results, it can find that Design I had a better identifying result than Design
10
11 308 II, especially based on GA method. For Design II, the classification rates of both Cu and Zn were
12
13
14 309 less than 90%, that was inferior to those of Cd and Pb.

15
16 310 As the improvement effect using CARS is a little better than those of using GA and SPA, an
17
18
19 311 in-depth analysis of characteristic variables was carried out based on CARS selection. SVM model
20
21 312 was established based on the characteristic variables selected by CARS algorithm. Figure 3 shows
22
23
24 313 the classification accuracy surface in a run of 10-fold CV, where the x - and y -axes represent $\log_2 C$
25
26 314 and $\log_2 \gamma$, respectively. Each mesh node in the (x, y) plane of the validation accuracy represents a
27
28
29 315 parameter combination, and the z -axis denotes the obtained classification accuracy value with
30
31 316 each parameter combination. Figure 3 shows that the classification rate with the optimal
32
33
34 317 parameters was over 90%. Figures 3 (a) to (e) show that the highest classification rate depended
35
36 318 on the optimization of the two parameters. In these figures, $C > 0$ and $\gamma < 0$ generally yielded high
37
38
39 319 classification rates. Meanwhile, the classification accuracy remained unchanged when the C value
40
41 320 exceeded a certain threshold. For Design I, with the optimum C and γ values, the variables
42
43
44 321 selected by CARS (33 for healthy samples vs. Cu-contaminated samples, 31 for healthy samples
45
46 322 vs. Zn-contaminated, 22 for healthy samples vs. Pb-contaminated samples, and 10 for healthy
47
48
49 323 samples vs. Cd-contaminated samples) were used as inputs of SVM to construct a classification
50
51 324 model. The model was validated by 10-fold CV using all samples. The classification rates during
52
53
54 325 training reached 90% (Figure 3). The corresponding prediction results are shown in Table 2. Table
55
56 326 2 shows that the effectiveness of the two-category classification (healthy samples vs. samples
57
58
59
60

1
2
3
4 327 contaminated by a single heavy metal) was highly satisfactory. This finding is evidenced by the
5
6 328 over 95% sensitivity, specificity and accuracy of the method. Compared with Design I, Design II
7
8
9 329 had slightly lower classification rate because greater complexity was required to identify more
10
11 330 sample varieties. Detailed classification results have been listed in Table 4, including
12
13 331 misclassifications of different samples, the prediction value for healthy samples was relatively
14
15
16 332 high (over 95%). However, several varieties of the heavy metal-contaminated samples were
17
18
19 333 recognized as healthy and most of them were Cu- and Zn-contaminated samples (four and five
20
21 334 samples, respectively). This observation may be attributed to the fact that unlike Pb and Cd, Cu
22
23
24 335 and Zn are essential to organisms. Little amounts of Cu and Zn may be insufficient to change the
25
26 336 protein and lipid structures of the samples. Slight changes in the IR spectra may increase false
27
28
29 337 classification rate. This phenomenon might be the reason for the lower classification rate for Cu
30
31 338 and Zn contaminations than for Pb and Cd contaminations.

339 **Analysis of Differences between Healthy and Contaminated Samples**

340 The spectral variables used in Designs I and II were selected using CARS which was utilized
341 10 times because a different model was constructed each time and the selected characteristic
342 spectral variables were slightly different. Figure 2 shows the characteristic spectral variables that
343 were selected five times or more by CARS. The IR spectrum of a biological system is derived
344 from the vibration of various functional groups; hence, the characteristic spectral information can
345 be attributed to the interactions between these groups. Such characteristic spectral information
346 reflects the specific characteristics of certain molecular structures. Figure 2 also shows that the
347 spectral variables selected by both Designs I and II are concentrated within the 624–1700 cm^{-1}
348 and 3000–3516 cm^{-1} regions.

1
2
3
4 349 An in-depth analysis shows that Design II selected more spectral variables than Design I. This
5
6 350 observation can be attributed to the fact that Design II requires more variables for model
7
8
9 351 construction because it must classify all five sample varieties. The variables selected by Design II
10
11 352 included almost all characteristic peaks within the full IR spectral range and distributed in two
12
13 353 main ranges, 1350–1540cm⁻¹ reflecting the COO-symmetric stretch (fatty acids and amino acids,
14
15
16 354 CH₂ bending: mainly lipids, and Amide II: N–H bending and C–N stretching of proteins) and
17
18
19 355 3500–3100cm⁻¹ reflecting Amides A and B. Besides, the 1680cm⁻¹ band that principally reflects
20
21 356 the Amide I and C=O stretching of proteins was also selected [11,14-15].
22

23
24 357 For Design I, the regions and numbers of the selected spectral variables differ between the
25
26 358 healthy samples and the samples contaminated by a single heavy metal. To differentiate healthy
27
28
29 359 from Cu-contaminated samples, the highest number of spectral variables (up to 33) was selected
30
31 360 using CARS method. These variables are concentrated at around 1500, 3100–3500, and 600cm⁻¹.
32
33
34 361 In total, 31 spectral variables were selected by CARS for the model that classifies healthy from
35
36 362 Zn-contaminated samples. These variables are concentrated at around 3000, 1600 and 1100cm⁻¹,
37
38
39 363 which are close to the bands of the variables for Cu contamination. Relatively fewer variables
40
41 364 were selected for Cd and Pb contaminations as compared with those for Cu and Zn contaminations.
42
43
44 365 For Pb contamination, 22 variables were selected, which mainly concentrate at the 2900–3100,
45
46 366 1500–1700, 1300 and 700–1100cm⁻¹ regions. For Cd contamination, only 10 variables were
47
48
49 367 selected, which are distributed at around 3300, 2800, 1600, 1100, 900 and 600cm⁻¹. The different
50
51 368 numbers of variables selected for the four heavy metal contaminants were due to the proposal that
52
53
54 369 Pb and Cd as highly toxic heavy metals can significantly change the structures of *Tegillarca*
55
56 370 *granosa* components (proteins, lipids and others). Therefore, less variable information is needed to
57
58
59
60

1
2
3
4 371 identify the samples contaminated by these two metals. Meanwhile, Cu and Zn are required at low
5
6 372 concentrations by most organisms and are unlikely to cause significant structural changes in
7
8 373 *Tegillarca granosa* because of their low toxicity. Hence, a more variable information is required to
9
10 374 identify samples contaminated by these heavy metals. From another perspective, the region at
11
12 375 around 1500 cm^{-1} was selected using CARS for the four sample varieties of heavy metal
13
14 376 contaminations. This region principally reflects the protein structure. This selection can be
15
16 377 explained by the mechanism of heavy metal poisoning. In other words, heavy metal poisoning
17
18 378 stimulates *Tegillarca granosa* to synthesize metallothionein. This compound causes peroxidation
19
20 379 reactions in the cell membrane and induces the synthesis of antioxidants, such as glutathione and
21
22 380 superoxide dismutase. These antioxidants lead to the generation of free radicals and H_2O_2 , which
23
24 381 are highly toxic to cells. The whole defense system collapses when the generation of H_2O_2 exceeds
25
26 382 the antioxidant capacity of an organism. Most of the structural changes caused by heavy metal
27
28 383 contamination are related to proteins [37-40]. This observation explains why the variables selected
29
30 384 for all four heavy metal contaminations include the region that reflects protein structure. A similar
31
32 385 explanation is also applicable to the variable selection using Design II. If more varieties of
33
34 386 contamination need to be classified, then more variables are required for model construction.
35
36
37
38
39
40
41
42
43
44

387 **Conclusion**

45
46
47 388 A rapid method based on IR spectroscopy and pattern recognition was proposed to
48
49 389 differentiate healthy samples from Cu-, Cd-, Zn-, or Pb-contaminated *Tegillarca granosa* samples.
50
51
52 390 The variables in the IR spectra that classified healthy and heavy metal-contaminated samples were
53
54 391 selected using CARS. These variables were used as inputs for SVM to construct a classification
55
56 392 model. The models that were classified between the healthy samples and the samples
57
58
59
60

1
2
3
4 393 contaminated by a single heavy metal showed accuracies of up to 95%. The method also classified
5
6 394 samples contaminated by different heavy metals with an accuracy of over 90%. This study
7
8
9 395 provides a new and convenient method for the rapid detection of heavy metal-contaminated
10
11 396 aquatic products.

14 397 **Acknowledgments**

17 398 The authors would like to acknowledge the financial support provided by National Natural
18
19
20 399 Science Foundation of China (NO.31201355).

23 400 **References**

- 26 401 [1] Institute for Nutrition and Food Safety of the Chinese Center for Disease Control and
27
28 402 Prevention, Food composition table: National representative value, Beijing, China, 1992, pp:
29
30 403 32-48.
- 33 404 [2] J. Xu, H. Zhou, X. Yan, X. Yan, C. Zhou, P. Zhu, and B. Ma, Effect of unialgal diets on the
35
36 405 composition of fatty acids and sterols in juvenile ark shell *tegillarca granosa* Linnaeus, *J.*
37
38 406 *Agric. Food Chem.*, 2012, **60**, 3973–3980.
- 41 407 [3] M. Barron, G. Stehly, and W. Hayton, Pharmacokinetic modeling in aquatic animals I .
42
43 408 models and concepts, *Aquat. Toxicol.*, 1990, **18**, 61-85.
- 46 409 [4] M. Soylak, and A. Aydin, Determination of some heavy metals in food and environmental
47
48 410 samples by flame atomic absorption spectrometry after coprecipitation, *Food Chem. Toxicol.*,
49
50 411 2011, **49**, 1242-1248.
- 54 412 [5] M. Tüzen, Determination of heavy metals in fish samples of the middle Black Sea (Turkey)
55
56 413 by graphite furnace atomic absorption spectrometry, *Food Chem.*, 2003, **80**, 119-123.

- 1
2
3
4 414 [6] R. Karoui, G. Downey, and C. Blecker, Mid-infrared spectroscopy coupled with
5
6 415 chemometrics: a tool for the analysis of intact food systems and the exploration of their
7
8 416 molecular structure structure-quality relationships - a review. *Chem. Rev.*, 2010, **110**,
9
10 417 6144-6168.
- 11
12
13 418 [7] X.J. Chen, D. Wu, X.C. Guan, B. Liu, G. Liu, M.C, Yan, H.L. Chen, Feasibility of Infrared
14
15 419 and Raman Spectroscopies for Identification of Juvenile Black Seabream (*Sparus*
16
17 420 *macrocephalus*) Intoxicated by Heavy Metals, *J. Agr. Food Chem.* 2013, 61, 12429-12435.
- 18
19 421 [8] D. Cozzolino, Recent trends on the use of infrared spectroscopy to trace and authenticate
20
21 422 natural and agricultural food products, *Appl. Spectrosc. Rev.*, 2012, **47**: 518-530.
- 22
23 423 [9] D. Wu, J. Chen, B. Lu, L. Xiong, Y. He, Y. Zhang, Application of near infrared spectroscopy
24
25 424 for the rapid determination of antioxidant activity of bamboo leaf extract. *Food Chem.*, 2012,
26
27 425 **135**, 2147-2156.
- 28
29 426 [10] D. Wu, P.C. Nie, Y. He, Y.D. Bao, Determination of calcium content in powdered milk using
30
31 427 near and mid-infrared spectroscopy with variable selection and chemometrics. *Food*
32
33 428 *Bioprocess Tech.*, **5**, 1402-1410.
- 34
35 429 [11] Z. Movasaghi, S. Rehman, and I. Rehman, Fourier transform infrared (FT-IR) spectroscopy
36
37 430 of biological tissues, *Appl. Spectrosc. Rev.*, 2008, **43**, 134–179.
- 38
39 431 [12] X. Lu, M. Webb, M. Talbott, J. Van Eenennaam, S. Doroshov and B. Rasco, A study of
40
41 432 biochemical parameters associated with ovarian atresia and quality of caviar in farmed white
42
43 433 sturgeon (*Acipenser transmontanus*) by Fourier Transform Infrared (FT-IR) Spectroscopy,
44
45 434 *Aquaculture*, 2011,**315**, 298-305.
- 46
47 435 [13] S. Akkas, M. Severcan, O. Yilmaz, and F. Severcan, Effects of lipoic acid supplementation on
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4 436 rat brain tissue: an FT-IR spectroscopic and neural network study, *Food Chem.*, 2007, **105**,
5
6 437 1281-1288.
7
8
9 438 [14] P. Palaniappan, and V. Vijayasundaram, Fourier transform infrared study of protein secondary
10
11 439 structural changes in the muscle of Labeorohita due to arsenic intoxication, *Food Chem.*
12
13 440 *Toxicol*, 2008, **46**:3534-3539.
14
15
16 441 [15] P. Palaniappan, and V. Vijayasundaram, The FT-IR study of the brain tissue of Labeorohita
17
18 442 due to arsenic intoxication. *Microchem. J.*, 2009, **91**,118-124.
19
20
21 443 [16] P. Palaniappan, and V. Vijayasundaram, The effect of arsenic exposure and the efficacy of
22
23 444 DMSA on the proteins and lipids of the gill tissues of Labeorohita. *Food Chem.Toxicol.*,
24
25 445 2009,**47**, 1752-1759.
26
27
28 446 [17] P. Palaniappan, and K. Pramod, FTIR study of the effect of nTiO₂ on the biochemical
29
30 447 constituents of gill tissues of Zebrafish (*Danio rerio*). *Food Chem. Toxicol.*, 2010, **48**:
31
32 448 2337-2343.
33
34
35 449 [18] H. Li, Y. Liang, Q. Xu, and D. Cao, Key wavelengths screening using competitive adaptive
36
37 450 reweighted sampling method for multivariate calibration. *Anal. Chim. Acta*, 2009, **648**, 77–
38
39 451 84.
40
41
42 452 [19] K. Zheng, Q. Li, J. Wang, J. Cheng, P. Cao, T. Sui, X. Wang and Y. Du, Stability competitive
43
44 453 adaptive reweighted sampling (SCARS) and its applications to multivariate calibration of
45
46 454 NIR spectra. *Chemometr. Intell. Lab.*, 2012, **112**, 48-54.
47
48
49 455 [20] W. Fan, Y. Shan, and G. Li, Application of competitive adaptive reweighted sampling method
50
51 456 to determine effective wavelengths for prediction of total acid of vinegar. *Food Anal.Method*,
52
53 457 2012, **5**,585-590.
54
55
56
57
58
59
60

- 1
2
3
4 458 [21] X. Zhang, W. Li, and B. Yin, Improvement of near infrared spectroscopic (NIRS) analysis of
5
6 459 caffeine in roasted Arabica coffee by variable selection method of stability competitive
7
8 460 adaptive reweighted sampling (SCARS). *Spectrochim. Acta A*, 2013, **114**, 350-356.
- 9
10
11 461 [22] D. Wu, Y. He, Potential of spectroscopic techniques and chemometric analysis for rapid
12
13 462 measurement of docosahexaenoic acid and eicosapentaenoic acid in algal oil. *Food Chem.*,
14
15 463 2014, **158**, 93-10.
- 16
17
18 464 [23] D. Wu, H. Shi, Y. He, X. Yu, Y. Bao, Potential of hyperspectral imaging and multivariate
19
20 465 analysis for rapid and non-invasive detection of gelatin adulteration in prawn. *J. Food Eng.*,
21
22 466 2013, **119**, 680-686.
- 23
24
25 467 [24] M.C.U. Araújo, T.C.B. Saldanha, R.K.H. Galvão, T. Yoneyama, H.C. Chame, V. Visani, The
26
27 468 successive projections algorithm for variable selection in spectroscopic multicomponent
28
29 469 analysis, *Chemometr. Intell. Lab.*, 2001, **57**, 65-73.
- 30
31
32 470 [25] R.K.H. Galvão, M.C.U. Araújo, W.D. Fragoso, E.C. Silva, G.E. José, S.F.C. Soares, H.M.
33
34 471 Paiva, A variable elimination method to improve the parsimony of MLR models using the
35
36 472 successive projections algorithm, *Chemometr. Intell. Lab.*, 2008, **92**, 83-91.
- 37
38
39 473 [26] K. Liu, X.J. Chen, L.M. Li, H.L. Chen, R.X. Ruan, W.B. Liu, A consensus successive
40
41 474 projections algorithm -multiple linear regression method for analyzing near infrared spectra,
42
43 475 *Anal. Chim. Acta*, 2015, doi:10.1016/j.aca.2014.12.033.
- 44
45
46 476 [27] D Wu, X.J. Chen, X.O., Zhu, X.C. Guan, G.C. Wu, Uninformative variable elimination for
47
48 477 improvement of successive projections algorithm on spectral multivariable selection with
49
50 478 different calibration algorithms for the rapid and non-destructive determination of protein
51
52 479 content in dried laver, *Ana. Method*, 2011, **3**, 1790-1796.
- 53
54
55
56
57
58
59
60

- 1
2
3
4 480 [28] R. Leardi, R. Boggia, M. Terrile, Genetic algorithm as a strategy for feature selection, J.
5
6 481 Chemometr. 1992, **6**, 267-281.
7
8
9 482 [29] D. Jouan-Rimbaud, D.L. Massart, R. Leardi, O.E. de Noord, Genetic algorithms as a tool for
10
11 483 wavelength selection in multivariate calibration, *Anal. Chim.*, 1995, **67**, 4295-4301.
12
13
14 484 [30] M. Pérez-Enciso, T. Tenenhaus, Prediction of clinical outcome with microarray data: a partial
15
16 485 least squares discriminant analysis (PLS-DA) approach, *Hum. Genet.*, 2003, **112**, 581-592.
17
18
19 486 [31] S. Wold, Pattern recognition by means of disjoint principal components models, *Pattern*
20
21 487 *Recogn.*, 1976, **8**, 127-139.
22
23
24 488 [32] E. Smidt, K. Meissi, M. Schwanninger, P. Lechner, Classification of waste materials using
25
26 489 Fourier transform infrared spectroscopy and soft independent modeling of class analogy,
27
28 490 *Waste Manage.*, 2008, **28**, 1699-1710.
29
30
31 491 [33] S. Balakrishnama, A. Ganapathiraju, Linear discriminant analysis - a brief tutorial, Institute
32
33 492 for Signal and information Processing, 1998, 1-8.
34
35
36 493 [34] C. Hsu, C. Chang, and C. Lin, A practical guide to support vector classification. In: Technical
37
38 494 report, Department of Computer Science and Information Engineering, National Taiwan
39
40 495 University, Taipei, 2003, available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.
41
42
43 496 [35] H. Li, Y. Liang, and Q. Xu, Support vector machines and its applications in chemistry.
44
45 497 *Chemometr. Intell. Lab.*, 2009, **95**, 188-198.
46
47
48 498 [36] H. Chen, C. Huang, and X. Yu, An efficient diagnosis system for detection of Parkinson's
49
50 499 disease using fuzzy k-nearest neighbor approach. *Expert. Syst. Appl.*, 2013, **40**, 263-271.
51
52
53 500 [37] G. Winston, and R. Di Giulio, Prooxidant and antioxidant mechanisms in aquatic organisms.
54
55 501 *Aquat Toxicol.*, 1991, **19**, 137-161.
56
57
58
59
60

- 1
2
3
4 502 [38] A. Viarengo, L. Canesi, M. Pertica, G. Poli, M. Moore and M. Orunesu, Heavy metal effects
5
6 503 on lipid peroxidation in the tissues of *mytilus-galloprovincialis* lam. *Comp. Biochem. Phys. C*,
7
8 504 1990, **97**, 37-42.
- 11 505 [39] R. Pipe, J. Coles, and F. Carissan, Copper induced immunomodulation in the marine mussel,
12
13 506 *Mytilusedulis*, *Aquat. Toxicol.*, 1999,**46**, 43-54.
- 16 507 [40] F. Regoli, G. Principato, Glutathione, glutathione-dependent and antioxidant enzymes in
17
18 508 mussel, *Mytilus-galloprovincialis*, exposed to metals under field and laboratory conditions:
19
20 509 implications for the use of biochemical biomarkers, *Aquat. Toxicol.*, 1995,**31**, 143-164.

21
22
23
24 510
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

511 **Figure caption**

512 Figure 1 Representative infrared spectra of healthy and heavy metal-contaminated samples of

513 *Tegillarca granosa*.

514 Figure 2 Variables selected using CARS for Designs I and II.

515 Figure 3 Training accuracy surface with parameters for the SVM model.

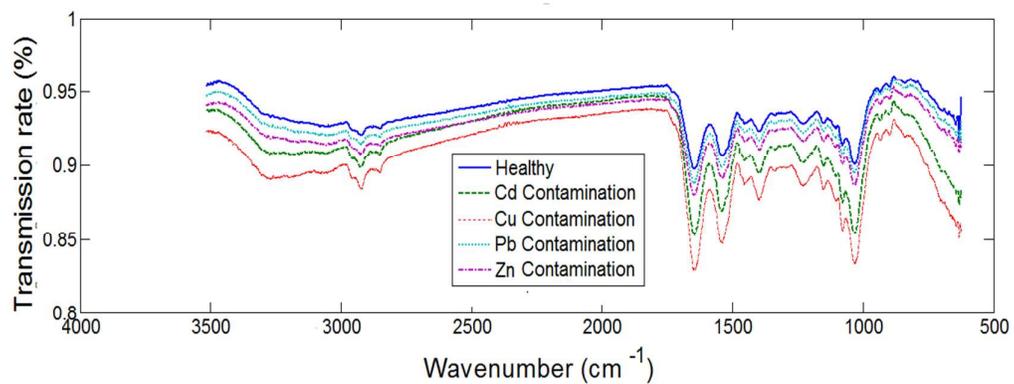


Figure 1 Representative infrared spectra of healthy and heavy metal-contaminated samples of

Tegillarca granosa.

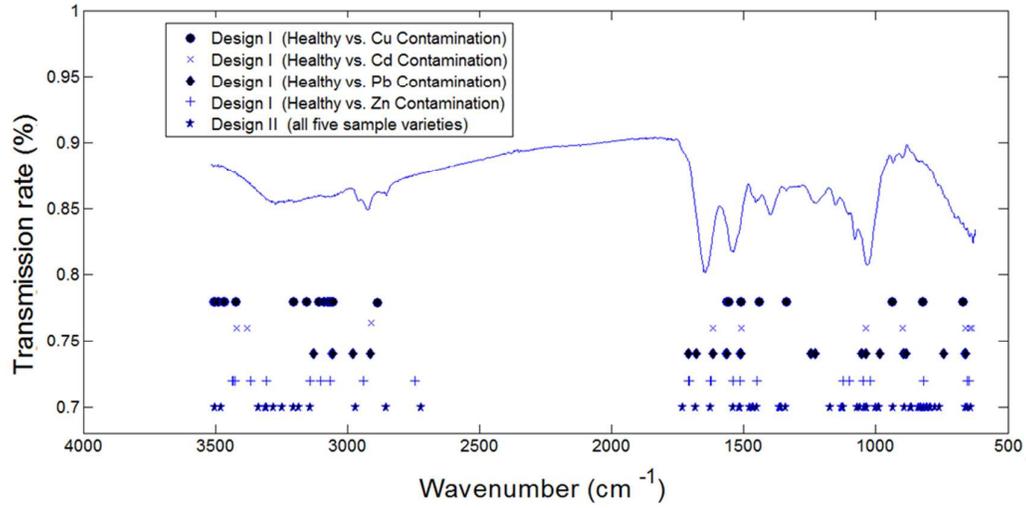


Figure 2 Spectral variables selected using CARS for Designs I and II.

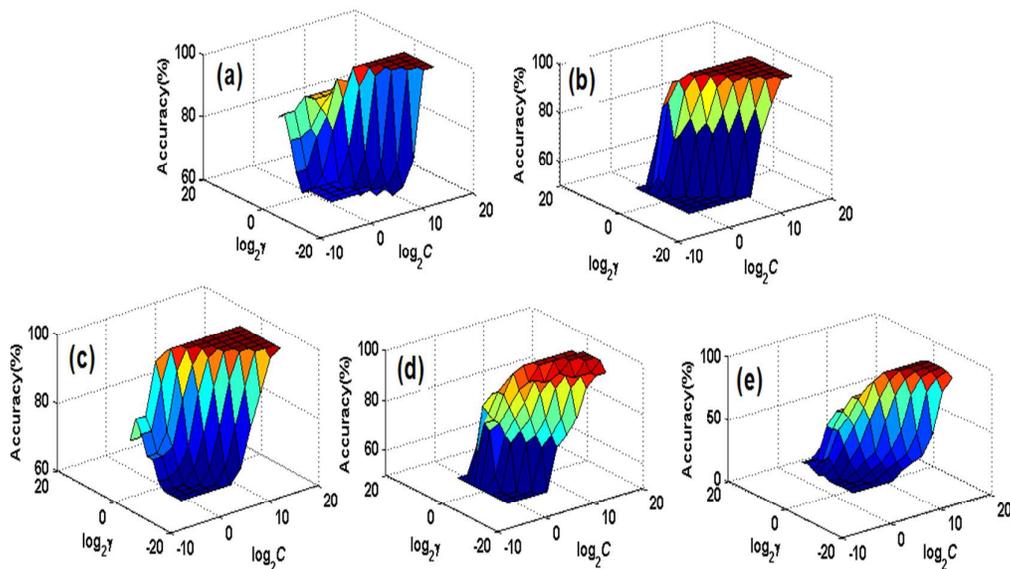


Figure 3 Training accuracy surface with parameters for the CARS-SVM model (a) Design I (Healthy vs. Cd), (b) Design I (Healthy vs. Cu), (c) Design I (Healthy vs. Pb), (d) Design I (Healthy vs. Zn), (e) Design II

Table 1 Heavy metal content statistical values of *Tegillarca granosa* samples

Sample	mean \pm SD(mg/kg)			
	Cd content	Cu content	Pb content	Zn content
Healthy	1.832 \pm 0.52	7.5 \pm 1.71	0.611 \pm 0.17	32.41 \pm 7.49
Contaminated by Cd	91.83 \pm 10.12	6.5 \pm 2.11	0.416 \pm 0.12	22.41 \pm 8.49
Contaminated by Cu	1.132 \pm 0.12	97.08 \pm 14.21	0.501 \pm 0.11	27.41 \pm 15.49
Contaminated by Pb	2.332 \pm 0.92	5.3 \pm 1.01	241.12 \pm 41.17	38.41 \pm 11.49
Contaminated by Zn	1.832 \pm 0.12	9.5 \pm 3.21	0.811 \pm 0.27	232.41 \pm 18.59

Table 2 Results of SVM model for Design I

	Method*	Training (%)			Prediction(%)		
		Sensitivity	Specificity	Accuracy	Sensitivity	Specificity	Accuracy
Healthy vs. Cd	CAR	94.8	97.7	96.3	97.0	95.0	96.3
	SPA	94.0	92.7	93.3	94.7	90.0	92.8
	GA	91.3	92.7	91.5	92.3	94.7	93.2
Healthy vs. Cu	CAR	98.7	95.3	96.0	95.6	94.3	95.3
	SPA	94.4	91.7	93.7	94.0	91.0	92.0
	GA	92.3	90.1	90.5	92.3	90.7	91.2
Healthy vs. Pb	CAR	95.3	98.0	96.7	95.3	97.3	96.1
	SPA	94.0	90.7	92.3	94.7	91.3	93.8
	GA	91.3	94.7	92.5	92.3	95.7	93.2
Healthy vs. Zn	CAR	98.7	95.7	97.7	97.0	94.4	95.9
	SPA	95.0	92.7	94.3	95.7	91.0	93.8
	GA	92.0	93.7	92.5	92.3	95.7	93.1

* Variable selection method

Table 3 Results of the CARs/SPA/GAs-SVM model for Design II

	Method*	Healthy	Cd	Cu	Pb	Zn
Training	CARs	95.1	96.0	93.9	96.6	92.8
	SPA	90.5	92.0	92.1	90.5	93.6
	GAs	94.7	91.6	92.9	91.6	92.1
Prediction	CARs	96.9	96.3	94.3	93.8	91.7
	SPA	92.7	93.3	91.7	92.7	93.3
	GAs	91.7	91.0	88.0	93.0	89.0

*Variable selection method

Table 4 Results of the CAR-SVM model for Design II

		Healthy	Contaminated by Cd	Contaminated by Cu	Contaminated by Pb	Contaminated by Zn
Training	Healthy	255	2	7	1	9
	Contaminated by Cd	2	262	3	2	2
	Contaminated by Cu	4	3	249	2	6
	Contaminated by Pb	3	4	1	259	3
	Contaminated by Zn	4	2	5	4	256
	Accuracy (%)	95.1	96.0	93.9	96.6	92.8
Prediction	Healthy	31	0	1	0	2
	Contaminated by Cd	0	26	1	1	0
	Contaminated by Cu	1	0	33	0	0
	Contaminated by Pb	0	1	0	30	0
	Contaminated by Zn	0	0	0	1	22
	Accuracy (%)	96.9	96.3	94.3	93.8	91.7