

## Proof-of-Principle of rTLC, an Open-Source Software Developed for Image Evaluation and Multivariate Analysis of Planar Chromatograms

Dimitri Fichou, Petar Ristivojevi#, and Gertrud Elisabeth Morlock

*Anal. Chem.*, **Just Accepted Manuscript** • DOI: 10.1021/acs.analchem.6b04017 • Publication Date (Web): 15 Nov 2016

Downloaded from <http://pubs.acs.org> on November 15, 2016

### Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a free service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are accessible to all readers and citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
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

# Proof-of-Principle of rTLC, an Open-Source Software

## Developed for Image Evaluation and Multivariate Analysis of Planar Chromatograms

Dimitri Fichou<sup>a</sup>, Petar Ristivojević<sup>a,b</sup> and Gertrud E. Morlock<sup>a,\*</sup>

<sup>a</sup>Chair of Food Science, Institute of Nutritional Science, and Interdisciplinary Research Center (IFZ), Justus Liebig University Giessen, Heinrich-Buff-Ring 26-32, 35392 Giessen, Germany

<sup>b</sup>On leave from the Innovation Center of the Faculty of Chemistry, University of Belgrade, PO Box 51, 11158 Belgrade, Serbia

\*Corresponding author. Tel.: +49-641-99-39141; fax: +49-641-99-39149; E-mail address: Gertrud.Morlock@uni-giessen.de (G. E. Morlock).

1  
2  
3 **ABSTRACT:** High-performance thin-layer chromatography (HPTLC) is an advantageous  
4 analytical technique for analysis of complex samples. Combined with multivariate data analysis,  
5  
6 it turns out to be a powerful tool for profiling of many samples in parallel. So far, chromatogram  
7  
8 analysis has been time-consuming and required the application of at least two software packages  
9  
10 to convert HPTLC chromatograms into a numerical data matrix. Hence, this study aimed to  
11  
12 develop a powerful, all in one open-source software for user-friendly image processing and  
13  
14 multivariate analysis of HPTLC chromatograms. Using the caret package for machine learning,  
15  
16 the software was set up in the R programming language with an HTML-user interface created by  
17  
18 the shiny package. The newly developed software, called rTLC, is deployed online and  
19  
20 instructions for direct use as web application, and in case required, for local installation are  
21  
22 available on GitHub. rTLC was created especially for routine use in planar chromatography. It  
23  
24 provides the necessary tools to guide the user in a fast protocol to the statistical data output (*e. g.*,  
25  
26 data extraction, preprocessing techniques, variable selection and data analysis). rTLC offers a  
27  
28 standardized procedure and informative visualization tools that allow the user to explore the data  
29  
30 in a reproducible and comprehensive way. As proof-of-principle of rTLC, German propolis  
31  
32 samples were analyzed using pattern recognition techniques, principal component analysis,  
33  
34 hierarchic cluster analysis and predictive techniques, such as random forest and support vector  
35  
36 machines.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

48 **KEYWORDS:** High-performance thin-layer chromatography; Multivariate analysis;  
49  
50 Chemometrics; Open-source software, R programming language; Caret package  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 **INTRODUCTION:** Natural extracts may contain thousands of individual compounds, and the  
4 majority of these are present in low concentrations down to the trace level. Though it is  
5 challenging, it is important to obtain reliable fingerprints that represent sound profiles of  
6 physiologically active compounds<sup>1</sup>. Its simplicity, cost-effective operation and the possibility of  
7 simultaneous analysis of up to 20 samples in parallel makes high-performance thin-layer  
8 chromatography (HPTLC) a technique of choice in herbal and food analysis<sup>2, 3</sup>. For evaluation,  
9 the HPTLC fingerprint of a complex sample is visually compared to that of a certified reference  
10 sample or to marker compounds being characteristic for the respective sample. The main  
11 disadvantage of such a manual pattern recognition technique and its visual comparison is its  
12 subjectivity, and it highly depends on the analyst's perception. Hence, hyphenation of HPTLC  
13 with high-sophisticated multivariate techniques provides objective fingerprints, mainly based on  
14 mathematical models<sup>4,5</sup>. As HPTLC chromatograms contain hundreds of pixels, this  
15 multidimensionality is used to extract a maximum of information out of the chromatograms<sup>4</sup>. For  
16 example, pattern recognition techniques can recognize chemical compound patterns, identify  
17 characteristic marker compounds as well as classify unknown samples according to their  
18 biological activity.

19  
20  
21  
22  
23  
24  
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  
Though increasing, there are still a limited number of research papers on the combination of  
HPTLC with multivariate data analysis. Most of these are based on the investigation of propolis,  
herbal samples, biopolymers and microalgae<sup>6-16</sup>. Although propolis is one of the most  
investigated honeybee product, the separation of its complex phenolic compound composition is  
still challenging analysts. After derivatization with Neu's reagent and detection at UV 366 nm,  
phenolic components showed differently colored bands. Such colorful HPTLC chromatograms  
are highly appropriate input data for evaluation by multivariate data analysis. There exists a wide

1  
2  
3 range of derivatization reagents with different specificity and capability of detection. The  
4  
5 resulting characteristically colored bands generate different profiles on the red, green and blue  
6  
7 (RGB) channels. Thus, derivatization reagents can substantially influence the separation  
8  
9 performance and data evaluation<sup>6</sup>.  
10

11  
12 Contrary to other chromatography techniques, such as high performance liquid chromatography  
13  
14 (HPLC) and gas chromatography (GC) which offer a direct export of data for further multivariate  
15  
16 analysis, images of HPTLC chromatograms must first be converted to a numerical data matrix.  
17  
18 Various software, toolboxes and algorithms have been applied for image processing and  
19  
20 multivariate analysis of HPTLC chromatograms so far (Table 1)<sup>6-16</sup>. Such packages lack in  
21  
22 domain-specific functionality, which results in a manual, lumbering and time-consuming  
23  
24 pipeline of the data handling. The user is forced to open, process and save the data through  
25  
26 different software packages and toolboxes to perform the analysis<sup>6-16</sup>.  
27  
28  
29

30  
31 For the first time, we describe and introduce rTLC in this study. It is a newly developed open-  
32  
33 source web application for image processing and multivariate analysis of HPTLC  
34  
35 chromatograms. The focus is laid on the different possibilities and advantages of the application,  
36  
37 such as a fast and simple image processing workflow and application of a range of chemometric  
38  
39 techniques suited for planar chromatography. One driving force for developing rTLC was to  
40  
41 provide users with a unique solution to analyze HPTLC data. The access to a simple and accurate  
42  
43 open-source web application, instead of purchasing a number of licenses, was another impetus.  
44  
45 Many useful features for the analysis of HPTLC data, such as line profile of target compounds,  
46  
47 band comparison, signal preprocessing as well as comma separated value (CSV) export for  
48  
49 analysis on other platforms were integrated. Pattern recognition techniques such as principal  
50  
51 component analysis (PCA), hierarchical cluster analysis (HCA) and heat map are applicable on  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 separate channels (RGB and gray scale) or in combination. Prediction techniques such as random  
4 forest (RF), linear discriminant analysis (LDA), support vector machine (SVM), partial least  
5 square (PLS) and classification and regression tree (CART) analysis were integrated as well. The  
6 increasing number of publications in the field of planar chromatography hyphenated with  
7 multivariate analysis motivated to redesign software and add many new tools. This makes rTLC  
8 suitable for a wide range of applications in herbal, food and environmental science.  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

## 20 EXPERIMENTAL SECTION

21  
22 **Set-Up of the Open-Source Web Application.** The rTLC application is written with the R  
23 programming language<sup>17</sup>. R is an open-source language and environment for statistical  
24 computing and graphics. A key feature of R lies in its community of sharing users, who  
25 contribute to the extension of the language via packages, allowing others to use their work. rTLC  
26 uses in particular the shiny package to create an HTML based user interface<sup>18</sup> and the caret  
27 package for machine learning<sup>19</sup>. This way, the application was deployed online and is directly  
28 accessible via a modern internet browser having internet connection. As it is a web application,  
29 the user needs not to install software. Direct use of rTLC ([http://shinyapps.ernaehrung.uni-  
31 giessen.de/rtlc](http://shinyapps.ernaehrung.uni-<br/>30 giessen.de/rtlc)), and in case required, instructions for local installation are available on GitHub:  
32 <https://github.com/DimitriF/rTLC-apps>.  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45

46 **Example Data Set.** A given sample set was used as proof-of-principle of the newly developed  
47 software. 106 samples of German propolis obtained from the Apicultural State Institute  
48 (Stuttgart, Germany) were analyzed in a previous study<sup>20, 21</sup>. The resulting 7 chromatograms in  
49 the JPEG format were manually labeled before the statistical analysis, leading to the assignment  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 of 37 blue-type and 69 orange-type samples of German propolis. The rTLC parameters set are  
4  
5 discussed subsequently.  
6  
7  
8  
9

## 10 **RESULTS AND DISCUSSIONS**

11  
12 rTLC, the newly developed open-source web application for image processing and multivariate  
13 analysis of HPTLC chromatograms, is introduced for the first time. The simple and streamlined  
14 workflow (Figure 1) provides the necessary tools to reproducibly guide the user in a fast protocol  
15 to the statistical data output. For regular cases, the evaluation of HPTLC chromatograms took  
16 only few minutes. The proof-of-principle was demonstrated via a German propolis data set,  
17 which was also made available as demonstration file in the rTLC software. Thus, the user is able  
18 to follow and reproduce the results reported below.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

29 **Data Input.** The user had to upload two files in rTLC to provide an appropriate data set for  
30 image evaluation and multivariate data analysis: (1) HPTLC chromatograms which contain the  
31 independent variables and (2) a batch file which contains the dependent variables about each  
32 sample on the plates, such as classes, botanical and geographical origin. Information on the  
33 experimental conditions is necessary to automatically extract each chromatogram from the  
34 HPTLC plate, *e.g.* the distances used during sample application and chromatography. rTLC  
35 supports the upload of commonly used image formats such as jpeg, tiff and png. The software  
36 computes the horizontal mean for each pixel of the chromatogram on the RGB channels as well  
37 as the gray scale, which is the mean of those three channels. At the end of this step, the data are  
38 in the form of a 3D array with samples as rows,  $R_F$  as columns and channels as layers (Figure 2).  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53 rTLC provides tools for line profiles of target compounds, comparison between tracks, pattern  
54 identification as well as identification of characteristic chemical and biological markers. The  
55  
56  
57  
58  
59  
60

1  
2  
3 profile comparison of RGB channels as well as gray scale helps to find similarities and  
4  
5  
6 dissimilarities between samples before and after signal preprocessing.

7  
8 **Data Preprocessing.** Recently, preprocessing methods used in HPTLC fingerprinting were  
9  
10 discussed<sup>22</sup>. Among others, the appearance of a non-homogeneous background after  
11  
12 derivatization, an increased noise level and band shifts are caused by variation in mobile and  
13  
14 vapor phase composition, humidity, temperature, operator handling and instrumental instability.  
15  
16 Thus, warping techniques are recommended to mitigate such experimental drawbacks<sup>4, 23</sup>. Two  
17  
18 peak alignment procedures were integrated into the rTLC software and are available to correct  
19  
20 inter- and intra-plate band shifts<sup>24</sup>: (1) parametric time warping and (2) dynamic time warping.  
21  
22 Further integrated options for data preprocessing such as denoising, normalization, and baseline  
23  
24 removal aimed at improving the quality of the data set. The software provides the Savitzky-  
25  
26 Golay and median filter, which are denoising/filtering methods commonly used in preprocessing  
27  
28 of HPTLC chromatograms<sup>4, 5</sup>. The baseline removal process was found to be mandatory in  
29  
30 almost all cases<sup>4</sup>, whereas good statistical models were also obtained without baseline  
31  
32 correction<sup>7</sup>. Hence, it is recommended to compare results with and without baseline correction.  
33  
34 Also, a normalization step is not mandatory and there is no consensus when it is obligatory -  
35  
36 sometimes it makes the results better, sometimes even worse. The preferred method of signal  
37  
38 normalization is the standard normal variate (SNV) method. Finally, rTLC provides auto-scaling  
39  
40 and mean centering to transform variables in the same unit<sup>5, 6, 24</sup>. The selection and need for  
41  
42 preprocessing tools depends on the project and may be chosen by the users to obtain ready-to-use  
43  
44 data for statistical analysis.  
45  
46  
47  
48  
49  
50  
51

52  
53 **Variable Selection.** HPTLC chromatograms provide a high number of variables for the given,  
54  
55 often limited number of available samples. There are several approaches regarding the nature of  
56  
57  
58  
59  
60



1  
2  
3 used variables for multivariate analysis. Important variables that contain information for the  
4  
5 aimed classification should be kept, whereas variables encoding the noise and/or with no  
6  
7 discriminating power should be removed<sup>25</sup>. For this purpose, rTLC provides options for careful  
8  
9 selection of variables for a specific channel or all channels together. The statistical analysis part  
10  
11 also informs on this selection, which can be optimized to keep the important information only.  
12  
13

14  
15 **Exploratory Statistics.** The user is only working with a data matrix, *i.e.* with samples as rows  
16  
17 and variables as columns; with this, it is possible to compute pattern recognition techniques such  
18  
19 as PCA, HCA and heatmap. For each of these techniques, informative visualization tools are  
20  
21 available that illustrate the data in various perspectives and allow the user to highlight patterns  
22  
23 by comparing the results with a chosen column of the batch file. For both, beginners and  
24  
25 experienced R-users, an editor is available and can be used for other types of techniques or  
26  
27 custom-made plots.  
28  
29

30  
31 **Predictive Statistics.** With the same matrix as mentioned before, this feature allows the user to  
32  
33 train a predictive model, used for the subsequent prediction of the properties of new samples.  
34  
35 There are two main techniques in predictive statistics, *i.e.* classification and regression; both are  
36  
37 available in the software. Before the training, the data set is split into training and test set to  
38  
39 produce a true validation set and avoid overfitting. The application uses the caret package<sup>19</sup> of  
40  
41 the R language to tune a model and choose the optimal parameters for a given algorithm. The  
42  
43 available predictive techniques are LDA, PCA (regression only), PLS (regression only), RF,  
44  
45 CART as well as SVM with linear and polynomial kernel.  
46  
47  
48  
49

50  
51 A model will be trained for each value of a grid, automatically created but editable, and the  
52  
53 parameters which give the best validation result will be kept for the final model. The choice of  
54  
55 the best set of parameters is made according to a cross-validation procedure; available  
56  
57  
58  
59  
60

1  
2  
3 procedures are k-fold cross validation, bootstrapping and leave-one-out-cross-validation. A  
4  
5 summary metric must be chosen to select the best model. For regression, the statistical  
6  
7 parameters can be expressed by root mean squared error (RMSE) or  $R^2$ . For classification,  
8  
9 accuracy, kappa, sensitivity or sensibility are available as summary metrics.

10  
11  
12 Different output tools are available to explore the result, such as confusion matrix of the test set,  
13  
14 prediction table and model summary. Also here, an editor is available to produce custom-made  
15  
16 plots. At the end of this step, a model file can be downloaded and used in other sessions to  
17  
18 predict the properties of new samples.  
19

20  
21  
22 **Proof-of-Principle of rTLC.** HPTLC chromatograms contain comprehensive information regard-  
23  
24 ding the polarity, chemical, and spectral properties of individual compounds in a sample. As a  
25  
26 case study, HPTLC chromatograms of German propolis samples were used to illustrate the  
27  
28 practical application of the rTLC software. The HPTLC chromatograms of propolis showed a  
29  
30 complex mixture of phenolic compounds, and thus, were highly appropriate input data to  
31  
32 demonstrate the performance and power of rTLC. Visual comparison of the respective HPTLC  
33  
34 chromatograms and RGB channels (Figure 3A) revealed a difference in the chemical  
35  
36 composition of the two types of German propolis. The blue type of propolis had several blue  
37  
38 bands at  $R_F$  values around 0.2, 0.3 and 0.6 (Figure 3B). The orange type of propolis showed a  
39  
40 rich phenolic profile and contained characteristic orange and yellow bands in the  $R_F$  range of 0.1-  
41  
42 0.5, and high fluorescent blue bands in the  $R_F$  range of 0.5-0.8 (Figure 3C). Next, two  
43  
44 unsupervised techniques (PCA and HCA) and two supervised techniques (RF and SVM) were  
45  
46 selected to illustrate the capabilities of multivariate analysis by rTLC. Parametric time warping  
47  
48 (aligned to the first sample), SNV and mean centering were used as preprocessing step.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 **Unsupervised Techniques.** Commonly used pattern recognition techniques<sup>5,6</sup> as PCA and HCA,  
4 are performed by rTLC in a fast and simple way. PCA was applied on the data set for the RGB  
5 channels as well as on the grayscale image. The variable of interest was class as color  
6 assignment (labelling and different symbols were not chosen). The blue channel (Figure 4A) and  
7 grayscale data (Figure 4D) with a  $R_F$  range of 0-1 as variable selection showed the best  
8 discrimination between the two sorts of German propolis samples and their statistical  
9 performances were discussed subsequently.

10  
11 In case of the blue channel data, PCA resulted in a five-component model, explaining 78.41% of  
12 the total variance. PC1 described 40.99%, while PC2 explained 15.34% of the total variance  
13 (Figure 4A). The most influential phenolic compounds were identified using the loading plots.  
14 For PC1, the compounds at  $R_F$  0.04, 0.38, 0.53, 0.66 and 0.98 had positive contributions while  
15 the compounds at  $R_F$  0.29, 0.58 and 0.77 had negative contributions (Figure 4B). For PC2, the  
16 compounds at  $R_F$  0.27, 0.52, 0.63 and 0.82 had positive contributions, while the compounds at  $R_F$   
17 0.06, 0.30, 0.36, 0.56 and 0.72 had negative contributions (Figure 4C).

18  
19 In the case of the grayscale image, the total variance explained by the first three PCs was 59.66%  
20 (PC1: 32.58%, PC2: 15.61%, and PC3: 11.45%) (Figure 4D). The discrimination between the  
21 two types of propolis samples is mainly driven by the first component. For PC1, positive  
22 influences were found at  $R_F$  0.06, 0.34, 0.39, 0.53 and 0.66 and negative ones at  $R_F$  0.28, 0.58  
23 and 0.79 (Figure 4E). For PC2, positive influences were observed at  $R_F$  0.04, 0.36, 0.57 and 0.73  
24 and negative ones at  $R_F$  0.27, 0.52, 0.65 and 0.84 (Figure 4F). Once those influential  $R_F$  values  
25 are known, the researcher can apply other analytical techniques or refer to the literature to  
26 identify such discriminatory compounds.

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

1  
2  
3 Cluster analysis is an often used classification technique. This algorithm performs a hierarchical  
4 cluster analysis using the distance between samples. At the beginning, each sample is assigned to  
5 its own cluster, iteratively, the closest clusters are joined together and the distances between  
6 clusters are recomputed, continuing until there is only one cluster. The simplest and most  
7 intuitive way to mathematically define the similarity between objects is based on the Euclidean  
8 distance. rTLC provides several routes to define the similarity between objects. According to the  
9 blue channel and grayscale data, there was a good discrimination between the orange- and blue-  
10 type propolis samples, which was in agreement with PCA<sup>5-7</sup>. For the blue channel, 'class bis' (x-  
11 labelling and color), Euclidean distance, ward method and a cluster number of 3 were chosen (2  
12 clusters for gray scale).

13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27 In the dendrogram of the blue channel data, there were three clusters (Figure 5A). The first  
28 cluster had a distance of 49 and was mainly composed of orange samples, whereas the second  
29 cluster had a distance of 25 and was dominated of blue-type propolis samples and the third  
30 cluster had a distance of 28 and consisted mainly of the latter samples. The dendrogram obtained  
31 for the grayscale data showed two clusters (Figure 5B). The first cluster contained almost all blue  
32 samples, while the second cluster consisted mostly of orange samples, the distance were  
33 respectively 57 and 62. The few blue-type propolis samples grouped into the orange-type cluster  
34 differed in their patterns compared to the other blue-type propolis samples. These samples can be  
35 considered as a mixture of both types of propolis due to the natural variation in the chemical  
36 composition. For such cases, it has to be proven that the variations in the experimental condition  
37 had been removed during the preprocessing step, as far as possible.

38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53 **Supervised Techniques.** In supervised techniques, a set of data describing objects of known  
54 features is used to construct a training set that is used to predict those features for new samples  
55  
56  
57  
58  
59  
60

1  
2  
3 then. Supervised techniques were applied in a wide range of chromatographic,  
4  
5 spectrophotometric and sensorial data, for quantification, fingerprinting, authentication and  
6  
7 detection of adulteration of food and herbal products<sup>25</sup>. The feature can be discrete, like the  
8  
9 geographical or botanical origin, or continuous like the concentration of a target molecule in the  
10  
11 investigated samples.  
12  
13

14  
15 As a first step of the supervised procedure, the data were split between training and test set.  
16  
17 Secondly, preprocessing techniques were applied on the training and test set. Note for  
18  
19 normalization, that the mean centering and standard deviation of the training set is used to  
20  
21 standardize the test set to avoid overfitting<sup>25</sup>. After the following variable selection, prediction  
22  
23 models were built using the training set for each row of the tuning grid and each step of the cross  
24  
25 validation procedure. Afterwards, the best parameters were selected and the final model was  
26  
27 trained with those parameters on the entire training set. Lastly, the reliability of the model was  
28  
29 evaluated using the test set. Two powerful supervised algorithms were selected to present this  
30  
31 feature: RF and SVM with linear kernel. Like for PCA and HCA, the following preprocessing  
32  
33 was used: parametric time warping, SNV and mean centering. In each case, the ratio of training  
34  
35 to test set was 3:1 and the cross validation method was 5-fold cross validation with total accuracy  
36  
37 as summary metric of choice for the selection of the best model. The outcome was studied for  
38  
39 each of the three channels and the grayscale image. In all cases, the prediction efficiency was  
40  
41 high and demonstrated the power of the technique to reproduce human decisions.  
42  
43  
44  
45  
46  
47

48  
49 Though RF has rarely been used as multivariate tool in food and herbal research so far, there are  
50  
51 several benefits that could make the RF algorithm an appropriate supervised tool in HPTLC  
52  
53 analysis: it can be used (1) when there are much more variables than observations, (2) for two- or  
54  
55 multi-classification and (3) for a good predictive performance, even when most of the predictive  
56  
57  
58  
59  
60

1  
2  
3 variables are noise, and thus a preselection of variables is not required. As another benefit, this  
4  
5 algorithm does not need standardization. The RF classifier needs optimization for two  
6  
7 parameters to generate a prediction model: the number of classification trees desired (ntree) and  
8  
9 the number of variables (mtry) which are used for tree growing in each tree. The accepted default  
10  
11 values for those two parameters are 500 for ntree, and sqrt (mall) for mtry, whereby mall is the  
12  
13 total number of variables in the original data set. The most important parameter, mtry, can be  
14  
15 optimized with the caret package, in contrast to the ntree parameter. This optimization led to  
16  
17 more accurate models<sup>26, 27, 28</sup>. By the way, the option PLS resulted in an equivalent outcome to  
18  
19 RF and SVM.  
20  
21  
22  
23

24 For all channels, the accuracy of classification of the training set was 100%. Those models were  
25  
26 clearly overfitted and this outcome must not be taken into account to judge a model. The  
27  
28 confusion matrix was obtained for each channel on the test set and during the cross validation  
29  
30 (Table 2 A). The green channel showed a good accuracy for cross validation and for the test set.  
31  
32 For the blue channel data and gray scale image, the comparison between cross validation and test  
33  
34 set showed more consistency, which was in accordance to PCA and HCA. Detailed statistical  
35  
36 parameters for the blue channel showed the performance of the model according to different  
37  
38 metrics (Table 2 B). The importance of the variables for the RF algorithm trained on the blue  
39  
40 channel is evident (Figure 6). In contrast to the variables highlighted in the loading plots of the  
41  
42 PCA, the model resulted in other variables to discriminate the two types of propolis.  
43  
44  
45  
46  
47

48 The SVM algorithm separates the classes by an optimal hyper plane that maximizes the distances  
49  
50 between classes by defining boundaries for the closest classes (support vectors) from the margins  
51  
52 of the class. This way, SVM minimizes the training error with regard to the separation of the  
53  
54 considered classes by using the least complex boundaries out of all possible ones. The optimal  
55  
56  
57  
58  
59  
60

1  
2  
3 hyper plane is obtained by an interactive algorithm that minimizes an error function that contains  
4 a parameter (penalty error) to control the complexity of the model and to avoid overfitting<sup>24, 25</sup>.  
5  
6 Even if the results for each channel and the grayscale image (Table 3 A) were comparable with  
7  
8 the RF results, this algorithm performed slightly worse in particular on the cross-validation data  
9  
10 set. The tuning step chose values of 0.25 for cost and 2 for gamma, except for the blue channel  
11  
12 where the optimum cost was 0.5. Detailed statistical parameters for the grayscale showed the  
13  
14 performance of the model according to different metrics (Table 3 B).  
15  
16  
17  
18  
19  
20  
21

## 22 CONCLUSIONS

23  
24 According to our knowledge, there was no dedicated all-in-one software for a streamlined image  
25  
26 evaluation and multivariate analysis of HPTLC chromatograms. The newly developed rTLC  
27  
28 application was designed as user-friendly open-source software to ease fingerprint comparisons.  
29  
30 New perspectives and conclusions on the data set are supported by a wide range of visualization  
31  
32 tools, owed to high plotting capabilities of the R software. A great step forward was achieved by  
33  
34 a substantial reduction of the analysis time. rTLC solved the supervised and unsupervised data  
35  
36 handling within few minutes, whereas the current practice needs several hours using at least two  
37  
38 different software packages. To the best of our knowledge, rTLC is the most concise tool  
39  
40 available for application of different pattern recognition and prediction techniques for HPTLC  
41  
42 chromatograms. On the one hand, the open-source asset of this application may attract users for  
43  
44 the powerful combination of HPTLC and multivariate analysis. On the other hand, it may  
45  
46 encourage the users to contribute to this technology through feedback, discussing ideas and  
47  
48 adding new functionalities to the software.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 **AUTHOR INFORMATION**  
4

5 **Corresponding Author**  
6

7  
8 \*E-mail: Gertrud.morlock@uni-giessen.de Tel. +49-641-99-39141. Fax +49-641-99-39149.  
9

10 **Notes**  
11

12 The authors declare no competing financial interest.  
13  
14

15 **ACKNOWLEDGMENTS**  
16

17 Thank is owed to the Ministry of Education, Science and Technological Development of the  
18 Republic of Serbia, grant No. 172017 for financial support of P.R.  
19  
20  
21  
22  
23  
24  
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



## References

- (1) Joshi D.D., *Herbal Drugs and Fingerprints*, Springer: New Dehli, India, 2012. DOI: 10.1007/978-81-322-0804-4.
- (2) Krüger S.; Morlock G.E. In *Instrumental Thin-Layer Chromatography, Handbooks in Separation Science*; Poole C., Ed; Elsevier Science: Amsterdam, 1st edition, 2014; pp. 409-429. DOI:10.1021/ac00073a721.
- (3) Morlock G.E.; Schwack W. *J. Chromatogr. A* **2010**, *1217*, 6600–6609. DOI: 10.1016/j.chroma.2010.04.058.
- (4) Komsta Ł. *Chromatogr. Res. Int.* **2012**, *Article ID 893246*, 1–5. DOI: 10.1155/2012/893246.
- (5) Milojković-Opsenica D.; Ristivojević P.; Andrić F.; Trifković J. *Chromatographia* **2013**, *76*, 1239–1247. DOI: 10.1007/s10337-013-2423-9.
- (6) Ristivojević P.; Andrić F.; Trifković J.; Vovk I.; Stanisavljević L.Ž.; Tešić Ž.L.; Milojković-Opsenica D. *J. Chemom.* **2014**, *28*, 301–310. DOI: 10.1002/cem.2592.
- (7) Morlock G.E.; Ristivojević P.; Chernetsova E.S. *J. Chromatogr. A* **2014**, *1328*, 104–112. DOI: 10.1016/j.chroma.2013.12.053.
- (8) Sârbu C.; Moț A.C. *Talanta* **2011**, *85*, 1112–1117. DOI: 10.1016/j.talanta.2011.05.030.
- (9) Sagi S.; Avula B.; Wang Y.H.; Zhao J.; Khan I.A. *J. Sep. Sci.* **2014**, *37*, 2797–2804. DOI: 10.1002/jssc.201400646.
- (10) Tian R.T.; Xie P.S.; Liu H.P. *J. Chromatogr. A* **2009**, *1216*, 2150–2155. DOI: 10.1016/j.chroma.2008.10.127.
- (11) Tang T.X.; Guo W.Y.; Xu Y.; Zhang S.M.; Xu X.J.; Wang D.M.; Zhao Z.M.; Zhu L.P.; Yang D.P. *Phytochem. Anal.* **2014**, *25*, 266–272. DOI: 10.1002/pca.2502.

- 1  
2  
3 (12) Ogegbo O.L.; Eyob S.; Parmar S.; Wang Z.T.; Bligh S.W.A. *Anal. Methods* **2012**, *4*, 2522–  
4  
5 2527. DOI: 10.1039/c2ay25373a.  
6  
7  
8 (13) Xie P.S.; Sun S.; Xu S.; Guo L. *J. Chromatogr. Sep. Tech.* **2014**, *5*, 249–258. DOI:  
9  
10 10.4172/2157-7064.1000249.  
11  
12 (14) Zarzycki P.K.; Zarzycka M.B.; Clifton V.L.; Adamski J.; Głód B.K. *J. Chromatogr. A* **2011**,  
13  
14 *1218*, 5694–5704. DOI: 10.1016/j.chroma.2011.06.065.  
15  
16 (15) Milojković Opsenica D.; Ristivojević P.; Trifković J.; Vovk I.; Lušić D.; Tešić Ž. *J.*  
17  
18 *Chromatogr. Sci.* **2016**, *54*, 1077–1083. DOI: 10.1093/chromsci/bmw024.  
19  
20 (16) Morlock G.E.; Ristivojević P. *Food Hydrocolloids*, in print. DOI:  
21  
22 10.1016/j.foodhyd.2016.10.005.  
23  
24 (17) R Core Team. *R: A language and environment for statistical computing*; R Foundation for  
25  
26 Statistical Computing: Vienna, Austria, 2016. <https://www.R-project.org/>.  
27  
28 (18) Chang W.; Cheng J.; Allaire J.; Xie Y.; McPherson J. *Shiny: Web application framework for*  
29  
30 *r*; 2016. <https://CRAN.R-project.org/package=shiny>.  
31  
32 (19) Kuhn M. *J. Stat. Softw.* **2008**, *28*, 1–26. DOI: 10.18637/jss.v028.i05.  
33  
34 (20) Kunz N.; Scholl I.; Schroeder A.; Morlock G.E. In *58th Annual convention of the*  
35  
36 *Association of the German Bee Research Institutes*, Berlin, Germany, March 29–31, 2011;  
37  
38 Poster P 21.  
39  
40 (21) Morlock G.E.; Scholl I.; Kunz N.; Schroeder A. *CAMAG Bibliogr. Service CBS* **2013**, *111*,  
41  
42 13-15.  
43  
44 (22) Ristivojević P.; Morlock G.E. *J. Planar Chromatogr.--Mod. TLC* **2016**, *29*, 310-317. DOI:  
45  
46 10.1556/1006.2016.29.4.10  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
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
- (23) Wong K.H.; Razmovski-Naumovski V.; Li K.M.; Li G.Q.; Chan K. *J. Pharm. Biomed. Anal.* **2014**, *95*, 11–19. DOI: 10.1016/j.jpba.2014.02.007.
- (24) Mohd K.S.; Azemin A.; Hamil M.S.R.; Bakar A.R.A.; Dharmaraj S.; Hamdan M.R.; Mohamad H.; Mat N.; Ismail Z. *Asian J. Pharm. Clin. Res.* **2014**, *7*, 110–116.
- (25) Berrueta L.A.; Alonso-Salces R.M.; Heberger K. *J. Chromatogr. A* **2007**, *1158*, 196–214. DOI: 10.1016/j.chroma.2007.05.024.
- (26) Breiman L. *Mach. Learn.* **2001**, *45*, 5–32. DOI: 10.1023/A:1010933404324.
- (27) Ai F.F.; Bin J.; Zhang Z.M.; Huang J.H.; Wang J.B.; Liang Y.Z.; Yu L.; Yang Z.Y. *Food Chem.* **2014**, *143*, 472–478. DOI: 10.1016/j.foodchem.2013.08.013.
- (28) Svetnik V.; Liaw A.; Tong C.; Culberson J.C.; Sheridan R.P.; Feuston B.P. *J. Chem. Inf. Comp. Sci.* **2003**, *43*, 1947–1958. DOI: 10.1021/ci034160g.
- (29) Chang C.-C.; Lin C.-J. *ACM Trans. Intell. Syst. Technol.* **2011**, *2*, 27:1–27:27. DOI: 10.1145/1961189.1961199.
- (30) Hamber V.; Irish H. *J. Chem. Educ.* **2007**, *84*, 842–847. DOI: 10.1021/ed084p842.

**Table 1**

Overview of publications related to HPTLC and multivariate analysis.

No.	Samples	Multivariate techniques	Software	Ref.
1	Herbs	PCA, Artificial neural network (ANN)	Matlab R2007 (MathWorks, Natick, MA, USA)	10
2	Propolis	PCA, HCA, Partial least square-discriminant analysis (PLS-DA)	Matlab R2011a (MathWorks, Natick, MA, USA), PLS toolbox version 6.2.1 (Eigenvector Research Incorporated, Manson, WA, USA) Image J1.48c version (Research Services Branch, National Institute of Mental Health, Bethesda, MD, USA.)	6
3	Herbs	PCA, PLS-DA Orthogonal PLS-DA (O-PLS DA)	SIMCA-P+ Version 12 (Umetrics AB, Umea, Sweden), VideoScan (CAMAG. Muttenz, Switzerland)	12
4	Propolis	PCA, HCA, LDA	Matlab R2011a (MathWorks, Natick, MA, USA), PLS toolbox version 6.2.1 (Eigenvector Research Incorporated, Manson, WA, USA) SPSS Version 21 (BM Corporation, Armonk, NY, USA), LIBSVM Version 3.16 <sup>29</sup>	7
5	Herbs	K-nearest neighbors Classification and regression tree (CART) Successive projection algorithm-linear discriminant analysis (SPA-LDA) PCA-discriminant analysis (PCA-DA) Support vector machine-discriminant analysis (SVM-DA), PLS-DA	Matlab R2012b (MathWorks, Natick, MA, USA) PLS toolbox version 7.3.1 (Eigenvector Research Incorporated, Manson, WA, USA) SPA toolbox 1.0(Homemade programs written in Matlab) Classification toolbox version 2.0 (Milano Chemometrics and QSAR Research Group, Milano, Italy)	21
6	Herbs	PCA	XLSTAT (Addinsoft, New York, NY, USA)	9
7	Herbs	PCA	Origin pro (OriginLab, Northampton, MA, USA)	13
8	Propolis	PCA, HCA	TLC Analyzer <sup>30</sup>	8
9	Propolis	Similarity analysis, HCA K-means clustering, ANN, SVM	Self-programmed software Xe2 IDE (Embarcadero, San Francisco, CA, USA), SPSS Version 21 (IBM Corporation, Armonk, NY, USA),LIBSVM Version 3.16 <sup>29</sup>	11
10	Biopolymers	PCA, HCA	Matlab R2011a (MathWorks, Natick, MA, USA), PLS toolbox version 6.2.1 (Eigenvector Research Incorporated, Manson, WA, USA) Image J1.48c version (Research Services Branch, National Institute of Mental Health, Bethesda, MD, USA.)	16

**Table 2**

RF algorithm model: confusion matrix for the test set of blue- and orange-type propolis samples as well as cross validation set for each channel (A) and detailed summary metrics for the blue channel on the three data sets (accuracy, sensitivity and specificity; B).

<b>A</b>			<b>Test set</b>			<b>Cross validation</b>		
<b>Channel</b>	<b>Optimum mtry</b>		<b>Blue-type</b>	<b>Orange-type</b>	<b>Accuracy</b>	<b>Blue-type</b>	<b>Orange-type</b>	<b>Accuracy</b>
Red	2	Blue	6	6	0.7857	14	11	0.8077
		Orange	0	16		4	49	
Green	15	Blue	9	3	0.8571	16	9	0.8590
		Orange	1	15		2	51	
Blue	2	Blue	9	3	0.8929	17	8	0.8333
		Orange	0	16		5	48	
Gray	2	Blue	9	3	0.8929	16	9	0.8590
		Orange	0	16		2	51	

<b>B</b>	<b>RF model parameters</b>	<b>Accuracy</b>	<b>Sensitivity</b>	<b>Specificity</b>
	Training set	1.0000	1.0000	1.0000
	Test set	0.8929	1.0000	0.8421
	Cross validation	0.8333	0.7727	0.8571

**Table 3**

SVM model with linear kernel: Confusion matrix for the test set of blue- and orange-type propolis samples and cross-validation set for each channel (A) and summary metrics for the grayscale image data on the three data sets (accuracy, sensitivity and specificity; B).

<b>A</b>	<b>Optimum</b>			<b>Test set</b>			<b>Cross validation</b>			
	<b>Channel</b>	<b>Cost</b>	<b>Gamma</b>	<b>Type</b>	<b>Blue-type</b>	<b>Orange-type</b>	<b>Accuracy</b>	<b>Blue-type</b>	<b>Orange-type</b>	<b>Accuracy</b>
	Red	0.25	2	Blue	6	6	0.7143	19	6	0.7692
				Orange	2	14		12	41	
	Green	0.25	2	Blue	8	4	0.8214	18	7	0.8077
				Orange	1	15		8	45	
	Blue	0.5	2	Blue	9	3	0.8214	17	8	0.7564
				Orange	2	14		11	42	
	Gray	0.25	2	Blue	8	4	0.8517	18	7	0.8590
				Orange	0	16		4	49	

<b>B</b>	<b>SVM model parameters</b>	<b>Accuracy</b>	<b>Sensitivity</b>	<b>Specificity</b>
	Training set	1.0000	1.0000	1.0000
	Test set	0.8571	1.0000	0.8000
	Cross validation	0.8590	0.8182	0.8750

1  
2  
3 **List of figures**  
4  
5

6 **Figure 1.** Workflow of the newly developed rTLC software performed within few minutes for regular cases.  
7  
8

9  
10 **Figure 2.** Processing of the experimental parameters for extraction of the HPTLC chromatograms to obtain the HPTLC densitograms.  
11  
12

13 **Figure 3.** RGB channels (A) and HPTLC chromatograms of the phenolic profiles of the blue-type (B) and orange-type (C) German  
14 propolis samples.  
15  
16  
17

18  
19 **Figure 4.** PC scores (A and D) and loading plots according to the blue channel (B and C) and grayscale image (E and F) evaluation.  
20  
21

22 **Figure 5.** Dendrograms for blue channel (A) and grayscale (B) image evaluation of the German propolis samples.  
23  
24  
25

26 **Figure 6.** Variable importance for the RF algorithm model trained with the blue channel in the discrimination of orange- and blue-  
27 type propolis samples (red: variables of PCA loading plots).  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49

Figure 1

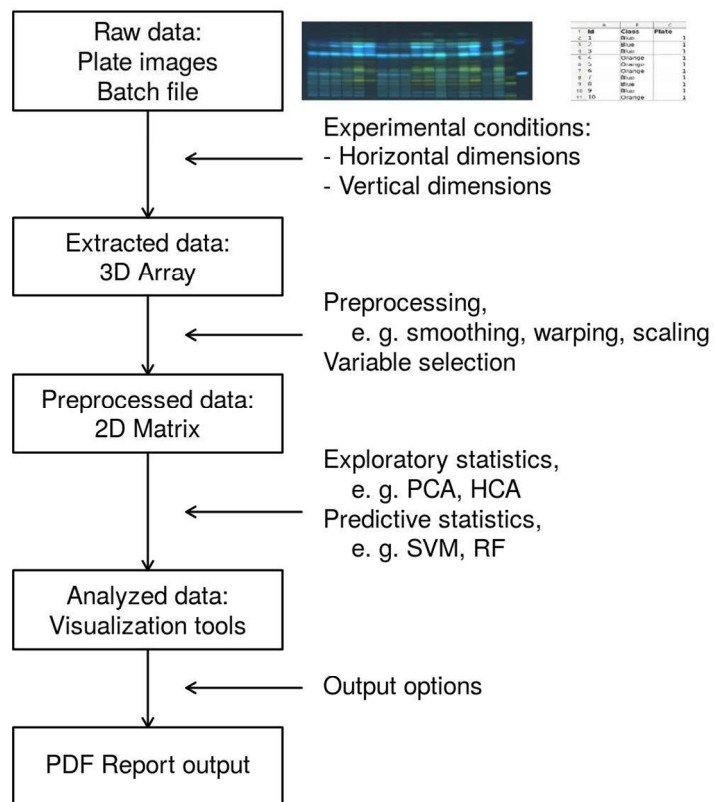




Figure 2

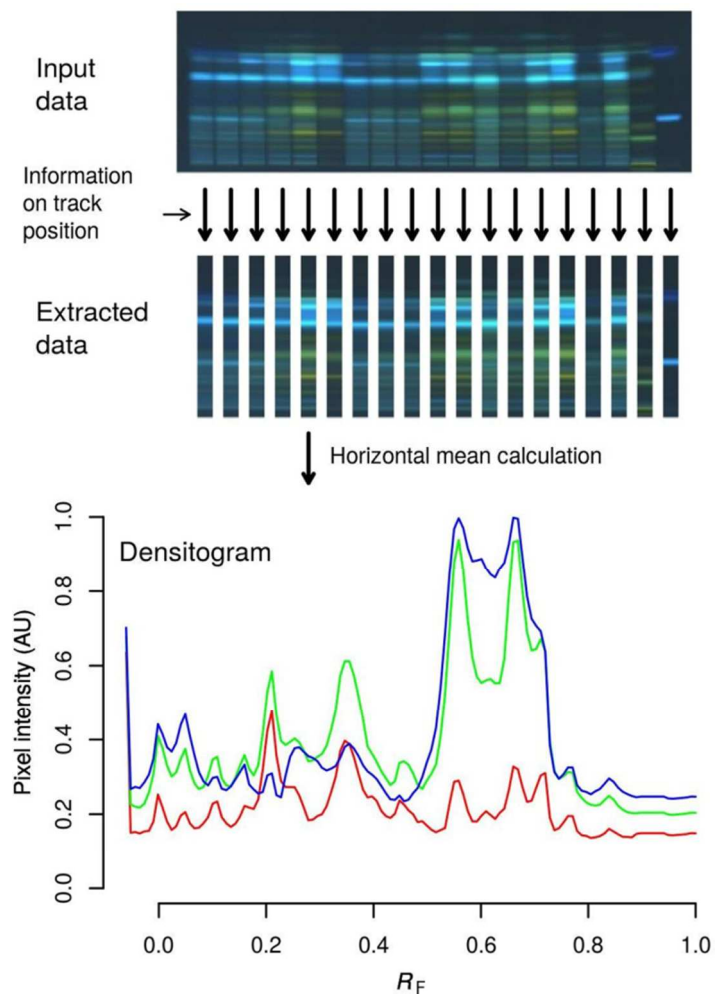


Figure 3

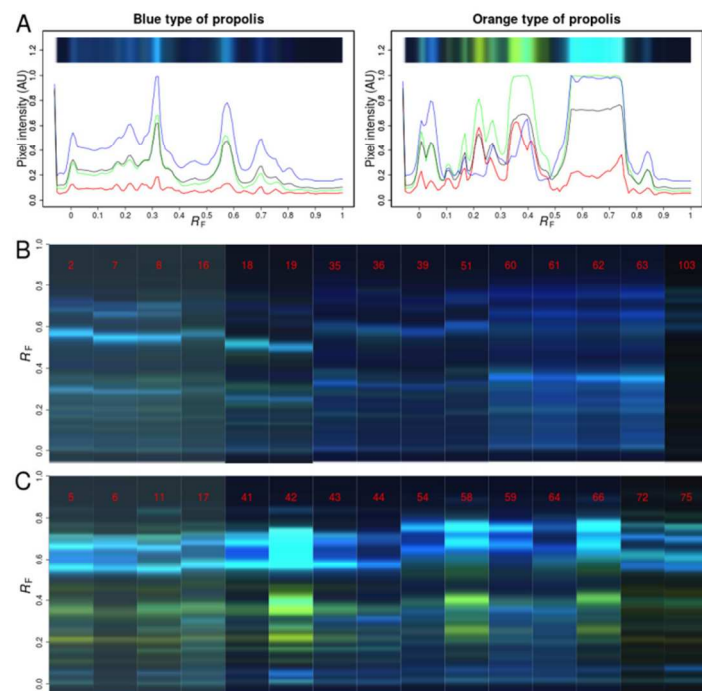


Figure 4

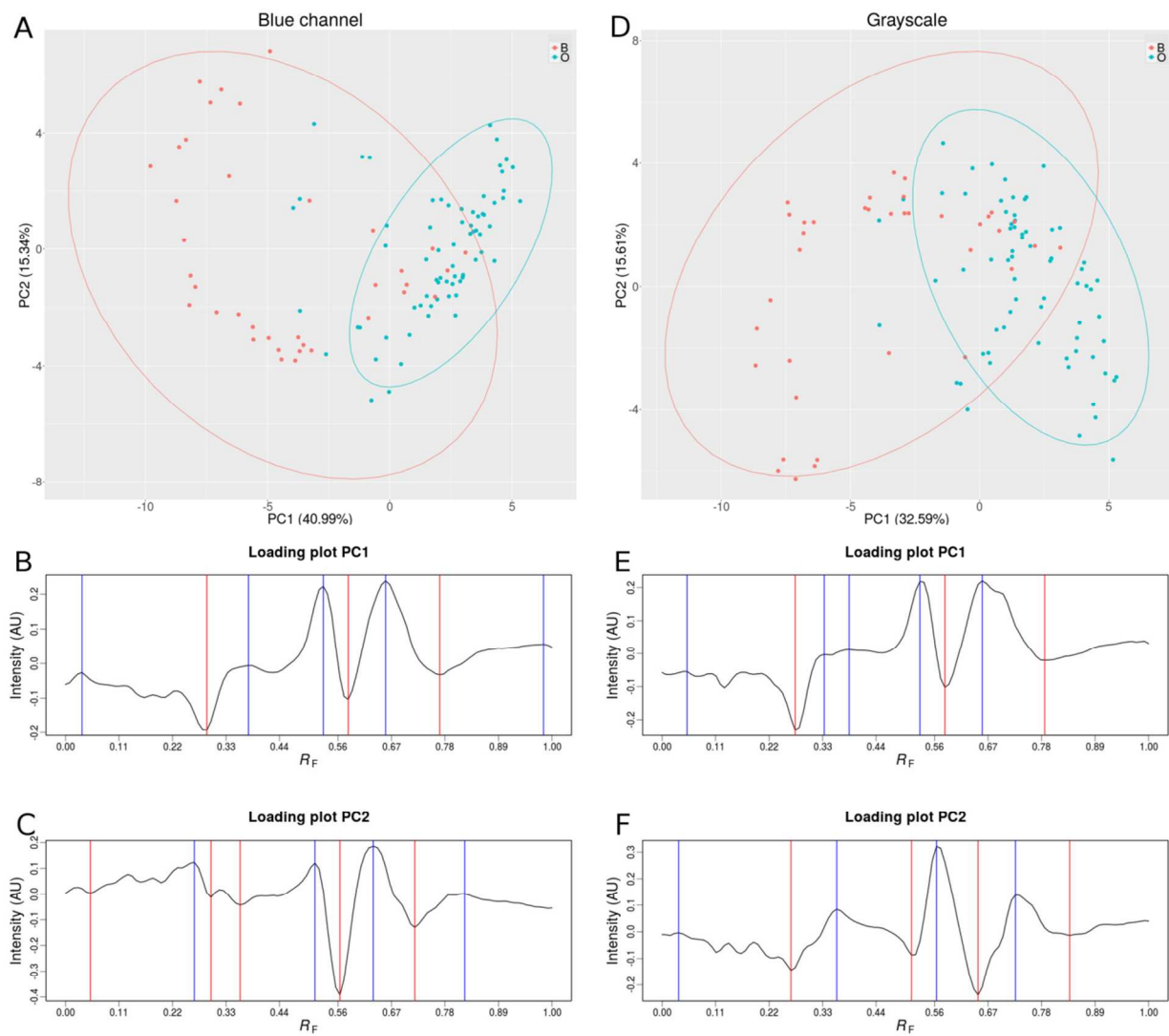


Figure 5

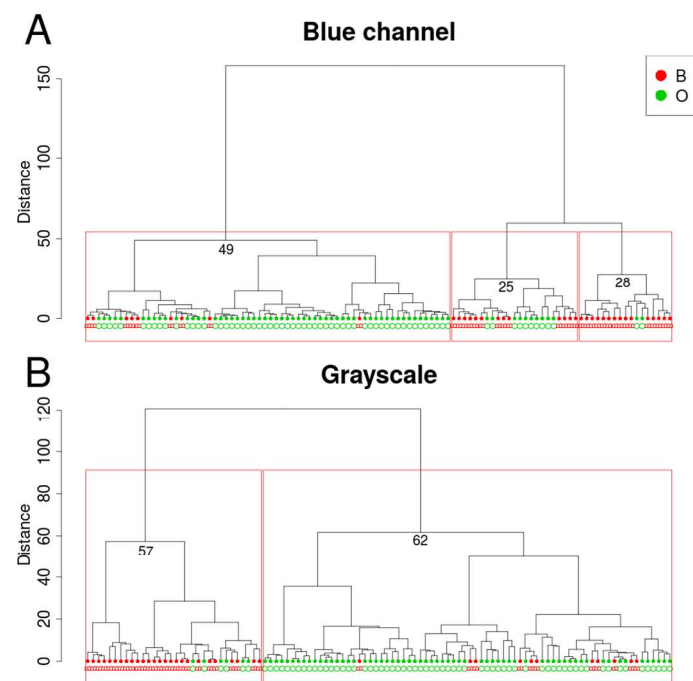
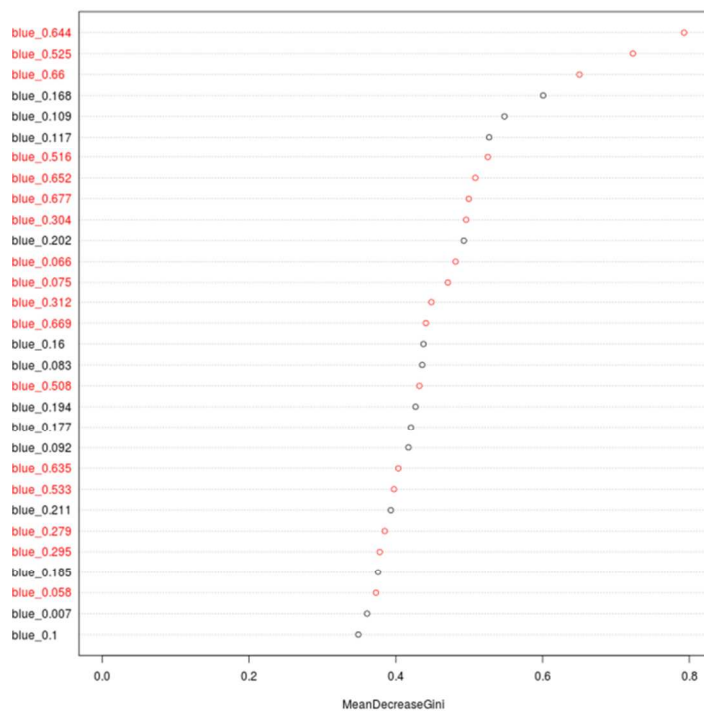


Figure 6



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
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

for TOC only

