



UNIVERSIDAD REGIONAL AMAZÓNICA IKIAM

Facultad de Ciencias de la Vida

Ingeniería en Biotecnología

Desarrollo de una librería estadística en R para el manejo de datos de cromatografía gaseosa bidimensional

Cristian Daniel Quiroz Moreno

20 de febrero de 2020, ciudad de Tena, Napo, Ecuador.

Declaración de derecho de autor, autenticidad y responsabilidad

Tena, 21 de febrero de 2020

Yo, Cristian Daniel Quiroz Moreno con documento de identidad N° 0604286443, declaro que los resultados obtenidos en la investigación que presento en este documento final, previo a la obtención del título Desarrollo de una librería estadística en R para el manejo de datos de cromatografía gaseosa bidimensional son absolutamente inéditos, originales, auténticos y personales.

En virtud de lo cual, el contenido, criterios, opiniones, resultados, análisis, interpretaciones, conclusiones, recomendaciones y todos los demás aspectos vertidos en la presente investigación son de mi autoría y de mi absoluta responsabilidad.

Por la favorable atención a la presente, suscribo de usted,

Atentamente,

Firma:



Cristian Daniel Quiroz Moreno
0604286443

Certificado de dirección de trabajo de integración curricular

Certifico que el trabajo de integración curricular titulado: “Desarrollo de una librería estadística en R para el manejo de datos de cromatografía gaseosa bidimensional”, en la modalidad de: artículo original, fue realizado por: Cristian Daniel Quiroz Moreno, bajo mi dirección.

El mismo ha sido revisado en su totalidad y analizado por la herramienta de verificación de similitud de contenido; por lo tanto, cumple con los requisitos teóricos, científicos, técnicos, metodológicos y legales establecidos por la Universidad Regional Amazónica Ikiam, para su entrega y defensa.

Tena, 20 de febrero de 2020

Firma:

Noroska Gabriela Sazalar Mogollón PhD
Directora del Laboratorio de Productos Naturales
C.I: 1757494453

Firma:

Guilherme L. Alexandrino PhD
Investigador Asociado en la Universidad de Copenhague

Agradecimientos

Agradezco a mi madre, Martha Moreno por todo su apoyo durante toda mi vida. Por la manera de enseñarme a vivir y las lecciones que he aprendido a lado de ella. Por el tesón indescriptible mostrado para poder criarnos a mí, a mi hermano. Le agradezco infinitamente por poder estar en esta etapa de mi vida.

Agradezco a mi hermano, Henry Quiroz, por todo su apoyo cuando decidí entrar a Ikiam. Por las veces que se sacrificó para que yo pudiese lograr un objetivo.

Agradezco a mi Padre, Rodrigo Quiroz, por sus enseñanzas y la superación que experimentó para ayudar a sus hijos. Además, a toda mi familia paterna por contribuir desde temprana edad para contribuir a mi personalidad.

Agradezco a todos mis familiares que me apoyaron en momentos difíciles como a Carmen Moreno, Elba Moreno; que dedicaron tiempo y esfuerzo para ayudar a mi familia en los momentos más duros. José Quiroz; que me ayudó en el proceso de profesionalización, entre muchos más que aportaron para poder lograr este objetivo.

Dedicatoria

Dedico este trabajo a grandes amigos, como Gabriel Gaona quien dedicó su tiempo para enseñarme una de las habilidades que me ha permitido destacar en mi campo de estudio. A Pablo Meneses, que en innumerables ocasiones intercedió para poder ocupar un lugar en el observatorio.

Dedico este trabajo a mi tutora, Gabriela Salazar, porque su mentoría se desarrolló en un ambiente de amistad y siempre estuvo para ayudarme, darme consejos sinceros ya sean académicos o personales, y no haber horario para conversar de lo que sea necesario.

Dedico a todos y cada uno de mis amigos que hicieron más amena la experiencia en esta etapa y apoyarme en buenas travesías.

Índice general

1. Introduction	3
2. Materials and methods.....	4
2.1. Methods	4
2.1.1. Biological material and its maintenance	4
2.1.2. Fungal Inoculation and Headspace Extraction	5
2.1.3. GCxGC-QMS.....	5
2.1.4. Tentative identification	5
2.2. Software implementation	6
2.2.1. Importing data	6
2.2.2. Visualization	10
2.2.3. Pre-processing.....	10
2.2.4. Multivariate analysis.....	12
2.3. Benchmarking.....	14
2.4. Application.....	14
3. Results and discussion	15
3.1. <i>Salmonella</i> dataset	15
3.2. <i>Penicillium</i> dataset.....	16
3.3. <i>Myrothecium</i> dataset	18
3.4. Computational efficiency	18
4. Conclusions	20

Índice de tablas

Table 1. Description of the main functions presented in the RGCxGC toolbox, classified by task aim. Each function has a short description, and the full argument list and default argument value. *The reference is included in the functions that are adapted from other, in cases when the function is built from scratch, the reference is left in blank.....8

Table 2. Description of the dataset analyzed with the proposed RGCxGC toolbox. While the dimensions of the single chromatogram are given by $A(I, J)$, the dimensions of the entire datasets are given by $A(I, J, K)$ where K represents the number of samples in each experiment.14

Índice de figuras

Figure 1. The proposed pipeline of non-targeted GCxGC-MS data analysis workflow in the RGCxGC toolbox. First, chromatograms are imported with the “read_chrom” function. Next, the user can pre-process them by smoothing, baseline correction, and peak alignment with the “wsmooth”, “baseline_corr” and “twod_cow” functions, respectively. Then, chromatograms from multiple cohorts are gathering in a single object, before to be subjected to multiway principal component analysis. On the other hand, the user is able to export all chromatograms with the “unfold_chrom”, in order to perform different statistical analysis. While the dashed lines enclose the functionalities implemented in the RGCxGC toolbox, the external parts show functionalities of external R toolboxes.....7.

Figure 2. The grid shows the results of the three datasets analyzed with the proposed toolbox; Salmonella (A, B, C, and D), Penicillium (E, F, G, and H) and Myrothecium (I, J, K, L). The first row shows the representative two-dimensional TIC chromatograms, the second row presents the MPCA scores, the third row displays the MPCA loadings, and the fourth row displays the PLS-DA scores. The results of unsupervised discriminant analysis by MPCA, of the Penicillium (F) and Myrothecium (G) datasets, present a clear separation between categories (edges in blue) in the two first PC’s, with an explained variance greater than 50%. On the other hand, the MPCA cannot explain the differences between the Salmonella (B) dataset categories. In the case of supervised classification with PLS-DA, the model can discriminate between categories for all datasets (D, H, and L).....17.

Figure 3. Benchmarking results of the main functionalities of the RGCxGC toolbox The head of each figure represents the function name that was tested. In order to create a similar situation of a big scale experiment, we tested the elapsed time for each function from 1 to 100 chromatograms samples in increments of 10 samples, except for the first step where the increment was 9 samples. All functions present a linear increment while the number of samples increases. Meanwhile “m_prcom” and “twod_cow” require 20 and 82 minutes to process 100 samples, the rest of the functions requires less than 5 minutes to process 100 samples..... 19.

Índice de anexos

Supplemental table. Tentative identification of detected metabolites by *Myrothecium sp.*
with spectral similarities >80%..... 36.

Resumen

La cromatografía de gases bidimensional completa (GC × GC) ofrece información química detallada sobre analitos volátiles y semivolátiles de muestras complejas. Sin embargo, la alta complejidad de la estructura de datos, fomenta el desarrollo de nuevas herramientas para un manejo y análisis de datos más eficiente. Aunque ya se han presentado algunas herramientas para superar este desafío, es necesario mejorar. En este manuscrito, presentamos una librería que contiene el flujo de trabajo para el procesamiento de datos GC×GC básico, que se puede utilizar tanto para el procesamiento previo de señales, como para el análisis de datos multivariados. Los algoritmos de preprocesamiento realizan el suavizado de la señal, la corrección de la línea de base y la alineación de los picos, mientras que el análisis multivariado se realiza a través del análisis de componentes principales (MPCA). El software es capaz de preparar los datos cromatográficos para otras aplicaciones con otras herramientas quimiométricas, por ejemplo: análisis de conglomerados, regresión, análisis discriminante, etc. El rendimiento de este nuevo software se probó en un conjunto de datos experimental interno y en otros dos conjuntos de datos publicados y disponibles en literatura.

Palabras clave: Metabolómica, análisis exploratorio, cromatografía bidimensional completa, caja de herramientas, procesamiento de señal.

Abstract

Comprehensive two-dimensional gas chromatography (GC×GC) offers detailed chemical information about volatile and semivolatile analytes from complex samples. However, the high complexity of the data structure encourages the development of new tools for a more efficient data handling and analysis. Although some tools have already been presented to overcome this challenge, there is still need for improvement. In this manuscript, we present a toolbox containing a pipeline for end-to-end basic GC×GC data processing which can be used for both, signal pre-processing and multivariate data analysis. The pre-processing algorithms perform signal smoothing, baseline correction, and peak alignment, while the multivariate analysis is done through Multiway Principal Component Analysis (MPCA). The software is capable to prepare the chromatographic data for further applications with other chemometric tools, e.g.: cluster analysis, regression, discriminant analysis, etc. The performance of this new software was tested on in-house experimental dataset and on two other published datasets.

Key words: metabolomics, exploratory analysis, comprehensive two-dimensional gas chromatography, toolbox, signal pr

Revista elegida

Nombre: *Microchemical Journal*

Página web: <https://www.journals.elsevier.com/microchemical-journal>

Referencia del artículo

Quiroz-Moreno, C., Furlan, M. F., Belinato, J. R., Augusto, F., Alexandrino, G. L., & Mogollón, N. G. (2020). RGCxGC toolbox: An R-package for data processing in comprehensive two-dimensional gas chromatography-mass spectrometry. *Microchemical Journal*, 156, 104830. <https://doi.org/10.1016/j.microc.2020.104830>

RGCxGC toolbox: an R-package for data processing in comprehensive two-dimensional gas chromatography-mass spectrometry

Cristian Quiroz-Moreno¹, Mayra Fontes Furlan², João Raul Belinato de Souza³, Fabio Augusto³, Guilherme L. Alexandrino⁴, and Noroska G.S. Mogollón^{1*}

¹Biomolecules discovery group, Universidad Regional Amazónica Ikiam, Tena-Ecuador

²Institute of Chemistry – University of Campinas (UNICAMP) and National Institute of Bioanalytical Science and Technology (INCTBio); CEP 6154, 13083-970 Campinas, – São Paulo, Brazil

³Institute of Chemistry, University of Campinas (UNICAMP), CEP 6154, 13083-970 Campinas, São Paulo, Brazil

⁴Department of Plant and Environmental Sciences, Faculty of Science, University of Copenhagen, Thorvaldsensvej 40, DK-1871

*Corresponding author:

Prof. Dr. Noroska G. S. Mogollón
Phone: +593 982310308

ABSTRACT: Comprehensive two-dimensional gas chromatography (GC×GC) offers detailed chemical information about volatile and semivolatile analytes from complex samples. However, the high complexity of the data structure encourages the development of new tools for a more efficient data handling and analysis. Although some tools have already been presented to overcome this challenge, there is still need for improvement. In this manuscript, we present a toolbox containing a pipeline for end-to-end basic GC×GC data processing which can be used for both, signal pre-processing and multivariate data analysis. The pre-processing algorithms perform signal smoothing, baseline correction, and peak alignment, while the multivariate analysis is done through Multiway Principal Component Analysis (MPCA). The software is capable of preparing the chromatographic data for further applications with other chemometric tools, e.g.: cluster analysis, regression, discriminant analysis, etc. The performance of this new software was tested on an in-house experimental dataset and on two other published datasets.

KEYWORDS: metabolomics, exploratory analysis, comprehensive two-dimensional gas chromatography, toolbox, signal processing

Declaration of interest: none

Abbreviations

2DCOW - Two-dimensional correlation optimized warping

CRAN – Comprehensive R archive network

DVB / CAR / PDMS - Divinylbenzene / Carboxene / Polymethylsiloxane

CUI – *Curvularia* sp.

DISCO – Distance and spectrum correlation optimization

GC – Gas chromatography

GC×GC - Comprehensive two-dimensional gas chromatography

GC×GC-MS - Comprehensive two-dimensional gas chromatography-Mass spectrometry

GUI – Graphical user interface

MPCA – Multiway principal component analysis

MS – Mass spectrometry

MYL – *Myrothecium* sp.

NetCDF - Network common data form

PC – Principal component

PCA – Principal component analysis

PDA – Potato-dextrose-agar

PLS - Partial least squares

PLS-DA - Partial least squares–discriminant analysis

PTW – Parametric time warping

SPME – Solid-phase microextraction

S/N – Singal to noise ratio

TIC - Total ion current

TOF – Time of flight

Q - Quadrupole

1. Introduction

Gas chromatography (GC), and particularly comprehensive two-dimensional gas chromatography, have come practically a mandatory technique for the analysis of the volatile and semivolatile compounds in matrices of high chemical complexity samples [1]. In GC×GC, the enhanced separation power is achieved by two capillary columns with preferably orthogonal separation capabilities, connected by the modulator. The modulator periodically concentrates a (coeluting) fraction of the eluate coming from the first column (first dimension, 1D) and next, reinjects this fraction as a narrower band into the second column (second dimension, 2D). Therefore, compounds that would coelute in conventional GC can be potentially separated in GC×GC system [2,3]. Because of the benefits provided by GC×GC, it has been widely applied in forensic [4], environmental [5], fuel [6], and metabolomics analysis [7].

Apart from the enhanced peak capacity of the GC×GC and the better elucidation of the chemical fingerprint from complex samples, isomers and homologous series are usually identified more easily in the two-dimensional chromatograms, e.g.: the roof-lite effect of hydrocarbons, and therefore the non-ambiguous identifications of unknowns is also improved [8,9]. However, GC×GC-MS can be coupled to multichannel detectors, such as mass spectrometers, which results in a large amount of data, e.g.: up to 2 GB raw files per sample in GC×GC-MS, and therefore efficient data handling tools are necessary [10–12].

Chemometrics use mathematical and statistical methods to analyze multivariate chemical data [13,14]. However, the efficient application of chemometrics in chromatography usually requires a previous pre-processing of the signals to reduce or mitigate undesirable artifacts, such as instrumental noise (e.g.: detectors' signal fluctuation and retention time shifts across multiple runs) and chemical noise (e.g.: column bleeding and peak saturation) [15]. The common pre-processing algorithms used in GC/GC×GC to correct for column bleeding and signal fluctuation are baseline correction and signal smoothing [15]. For GC×GC, the algorithms for correction of retention time shifts across samples should handle shifts in both ¹D and ²D, such as the two-dimension correlation optimized warping (2D-COW) [16], parametric time warping (DTW) [17], distance and spectrum correlation optimization (DISCO and DISCO2) [18,19], graph-based multiple alignment (BIPACE 2D) [20].

Data handling in GC×GC-MS, especially the pre-processing step, is usually performed by commercial software, and a minor percentage is done using open source tools [21,22]. For instance, over 85% of data analysis has been performed with the vendor or commercial software in metabolomics research, as manually calculated from the supplemental material from [23]. Although the open-source toolbox Guineu [24] and R2DGC [24] have been developed for data handling and pre-processing of GC×GC using single-and/or multichannel detection, there is still a need for extended workflow pipeline that also performs data analysis using chemometrics. In this stream, the RMet toolbox has been proposed as an end-to-end pipeline to process GC/GC×GC-MS data [25]. One of the main advantages of RMet is that it presents a Graphical User Interface, which makes the analysis easier

for inexperienced users. In contrast with command line toolboxes, GUI based ones do not allow an easy way to run multiple sets of parameter combinations until convergence is reached, as usually needed in order to tune parameters. Moreover, in the RMet toolbox, discriminant analysis is performed only by building a partial least squares-discriminant analysis model. Supervised models, as PLS-DA, have to be built carefully, since one of the main disadvantages is model overfitting [26].

Data analysis in GC/GC×GC has been commonly performed by pixel-based level, or by peak picking approaches [11,27,28]. The pixel-based approach performs the data analysis directly in the two-dimensional chromatogram which can be considered as a pixel image, obtained at the retention times in ¹D and ²D, respectively [29]. Conversely, the peak picking approach picks, integrates, and then organizes the individual peaks of the chromatograms in a peak table where the variances of the areas can be analyzed between multiple samples. There are advantages and drawbacks for both, pixel-based and peak picking approaches, and the decision about which approach to choose depends on the goals of the study. For instance, data analysis based on peak picking handles significantly fewer variables than in the pixel-based approach, which is an advantage for data analysis using the conventional univariate statistic. However, the quality of the data depends on the efficient selection and integration of the peaks, which can be problematic for highly coeluting peaks. Furthermore, the two-dimensional structure of the chromatograms is not seen straightforwardly as with a pixel-based approach is. Therefore, interpreting the chromatographic differences between samples can be less intuitive than the pixel-based approach.

In this work, we present the toolbox RGC×GC which was developed for data processing in GC×GC-MS, based on the open-source R environment. This toolbox contains an end-to-end pipeline for the most common processing techniques that are required for GC×GC, such as baseline correction, signal smoothing, and two-dimensional peak alignment. Moreover, the pixel-based analysis can be performed using the multiway principal component analysis (MPCA). On the other hand, the data can be exported in a more compatible format to be used with external toolboxes for chemometrics. In order to validate the RGC×GC toolbox, the performance of this open-source is demonstrated with three datasets in total. Two datasets were retrieved from literature, and one was in-house created.

2. Materials and methods

2.1. Methods

2.1.1. Biological material and its maintenance

The *Myrothcium sp* strains were obtained from "Collection of Bahia Microorganisms " (CCMB) with 0069/06 identifier. This strain was isolated by researches involved in the SisBiota project "Saprobic fungi bioprospection in the northeast of PPCIO/semi-arid region for the control of infectious diseases in plants: resistance induction" led by Sérgio Florentino Pscholati PhD. The strain was kept in Petri dishes containing 20 mL of carrot-maize-agar (CMA) culture media at (25,0 ±

1,0) °C in a growth chamber (Eletrolab, model EL202) with 12 h of photo-period.

2.1.2. Fungal Inoculation and Headspace Extraction

PDA culture media (25 mL) was placed in 50 mL polypropylene centrifuge tubes using angulation (elevation of 1,5 cm). The tube cap was modified with a 15 mm diameter hole and PTFE septa held by the aluminum ring.

Inoculation was made from Petri dishes with the fully grown fungal cells, and sterile distilled water was used to wash the surface of the plate. The plate was scraped with a sterile glass handle to obtain the spore suspension. The suspension was liquated and the concentration of 2.4×10^5 spores / mL was determined using a Neubauer chamber and an optical microscope. Then, 50 µL of the suspension was inoculated in a flow chamber into the tubes containing the culture medium. Tubes were kept at $(25,0 \pm 1,0)$ °C in a growth chamber (Eletrolab, model EL202) with 12 h of photoperiod.

A solid-phase microextraction (SPME) was performed from the third day of culture to the seventh day of culture. The SPME assay containing a DVB / CAR / PDMS (Divinylbenzene / Carboxene / Polymethylsiloxane 50/30 mm, Supelco) fiber was placed into the tube headspace for 35 min at $(25,0 \pm 1,0)$ °C.

2.1.3. GCxGC-QMS

A set of columns consisting of HP-5MS 30 m \times 0,25 mm \times 0,25 µm (Supelco) connected to a Supelcowax 1 m \times 0.10 mm \times 0.10 µm (Supelco) with a 1 m \times 0.25 mm deactivated silica capillary being used as the loop. The modulation period was set to 5.0 s. For GCxGC-QMS experiments were used a temperature program with the following parameters, 60 °C - 165 °C @3 °C/min; 165 °C - 260°C @20 °C/min; 260 °C (5 min); flow rate 0,6 mL/min (Helium 5.0 carrier gas); splitless injection mode, ion source temperature 200°C, interface temperature 260°C; voltage 0,9 kV; mass range 50-380 m/z; acquisition rate 25Hz and electron ionization (70 eV). For GCxGC-QMS data acquisition, GCMSsolution version 5.3 software (Shimadzu, Tokyo, Japan) and GCImage version 2.0 software (Zoex - Houston, TX, USA) were used for the analysis of two-dimensional chromatograms.

2.1.4. Tentative identification

For tentative identification, a two-dimensional chromatographic run was performed before the experiment with standards of homologs n-alkanes (C₈-C₂₀, Sigma Aldrich). 2 µL of the standard solution was transferred to a vial and the SPME fiber was exposed during 15 min after chromatographic desorption. The NIST 2008 spectra library (NIST – Gaithersburg, MD, EUA) was used (considering 80% of similarity) and comparison was made with the van den Don and Kratz retention index [30].

2.2. Software implementation

The basic workflow of the RGCxGC package is composed of three main steps; data importing, pre-processing and multivariate analysis. First, the raw Network Common Data Form (NetCDF) chromatogram is imported with the “read_chrom” function, in which the user needs to set the modulation time in which the GCxGC data was acquired. Next, you can perform smoothing and baseline correction using the function “wsmooth” and “baseline_corr”, respectively. Then, peak alignment from a single sample can be done using the “twod_cow” function, based on the two-dimensional correlation optimized warping (2DCOW) algorithm. Alternatively, multiple sample alignments can be performed with the “batch_2DCOW” routine, where the first chromatogram will be considered as the reference while aligning the remaining chromatograms. After pre-processing, MPCA can be performed on the dataset using the “m_prcomp” function, which provides the scores and loadings matrices and the summary with the explained and cumulative variance per Principal Component (PC). In the case of the loading matrix, they can be plotted using “plot_loading” and retrieved with the “scores” functions, while the “print” function you can access to the MPCA summary. The toolbox pipeline is summarized in Fig. 1. Library methods, and their arguments and comments, are summarized in Tab 1.

2.2.1. Importing data

This initial step is an adaptation of the Skov routine [31], this procedure is about how to extract and handle chromatographic signals. In the importing function, an option to import specific retention times ranges in both dimensions were included (*see x_cut and y_cut*) see Tab. 1. First, the chromatogram has to be exported from vendor software into a NetCDF file. This file extension is commonly used in scientific approaches [32]. The exported chromatogram contains the retention time and the Total Ion Current (TIC). Due to NetCDF architecture, data is stored into one-dimensional arrays, therefore, the signal is accessed through the arrays named scan_acquisition_time and total_intensity. Since the acquired data by both mass analyzers, time of flight and quadrupole, converge in a NetCDF file, they can be imported and analyzed with the proposed toolbox. Thus, the retention time is divided by sixty to transform the signal from seconds to minutes. Before the one-dimensional array is folded into a more familiar two-dimensional chromatogram, the sampling rate is evaluated to be homogenous. In other words, the software ensures that the entire run has the same sampling rate since non-integer sampling rates leads to unpaired data points over the chromatographic run.

Once the one-dimensional vector is stored in memory, the routine proceeds to fold it into a two-dimensional chromatogram. Each modulation period creates a matrix $C_{(I, J)}$, where I is the mass spectra scans, acquired in a given modulation and J is the modulation index. The number of mass spectra scans is related to the sampling rate (Hz) of the mass analyzer. Therefore, in order to calculate the total number of modulations, the total number of scans is divided by the number of scans per modulation.

Finally, once the number of scans per modulation (I) and the total number of modulations (J) are calculated, the one-dimensional array is folded into a two-dimensional matrix. Usually, a chromatographic run contains incomplete modulation scans at the end of the run, these scans are removed previous to create the two-dimensional chromatogram while printing a message.

In GC, solvent effect and column bleeding are almost unavoidable [1]. Therefore, in order to remove parts of the chromatogram that do not contain significant information, the user can provide the retention time range that they would like to keep. Another example is large chromatographic runs, where a cleaning ramp at high temperatures is included at the end of the run, producing large values at ¹D. Here, the user can avoid importing the last part of the chromatogram by changing the “x_cut” argument.

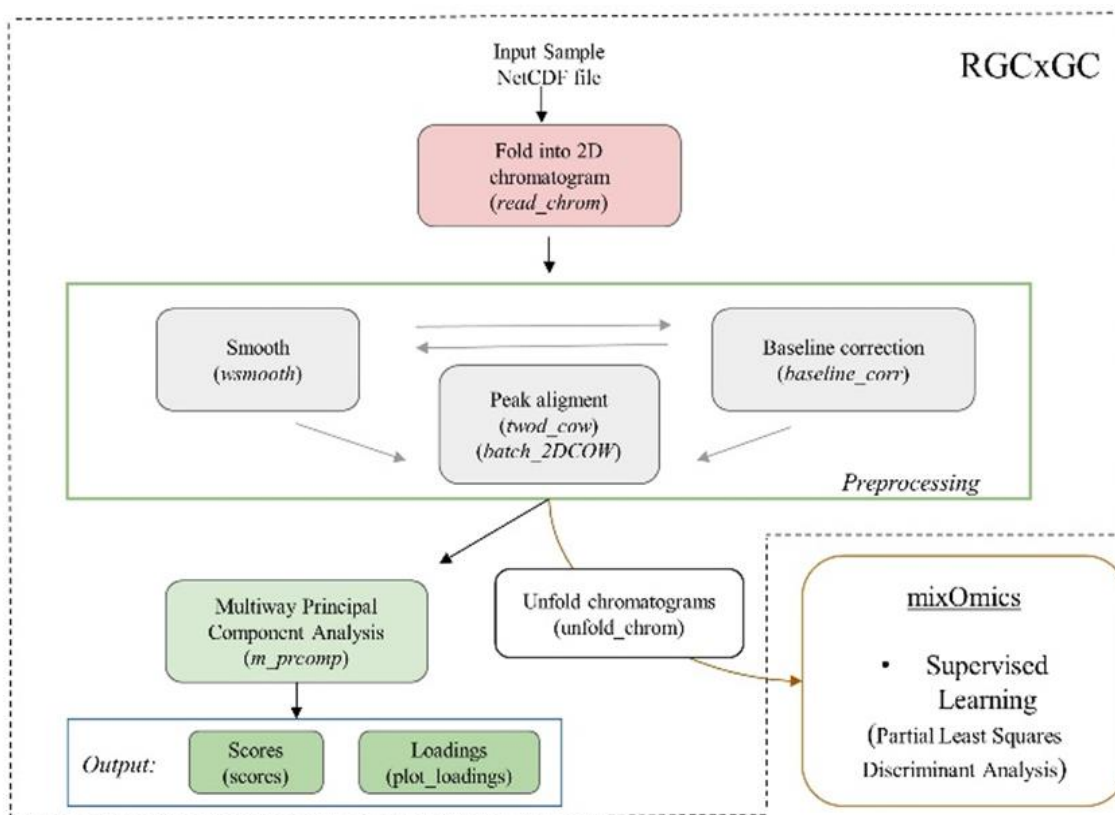


Figure 1. The proposed pipeline of non-targeted GCxGC-MS data analysis workflow in the RGCxGC toolbox. First, chromatograms are imported with the “read_chrom” function. Next, the user can pre-process them by smoothing, baseline correction, and peak alignment with the “wsmooth”, “baseline_corr” and “twod_cow” functions, respectively. Then, chromatograms from multiple cohorts are gathering in a single object, before to be subjected to multiway principal component analysis. On the other hand, the user is able to export all chromatograms with the “unfold_chrom”, in order to perform different statistical analysis. While the dashed lines enclose the functionalities implemented in the RGCxGC toolbox, the external parts show functionalities of external R toolboxes

Table 1. Description of the main functions presented in the RGCxGC toolbox, classified by task aim. Each function has a short description, and the full argument list and default argument value. *The reference is included in the functions that are adapted from other, in cases when the function is built from scratch, the reference is left in blank.

Aim	Task	Function name	Function description	Arguments	Argument description	Reference*
Importing data and visualization	Read raw chromatogram	read_chrom	This function reads the netCDF file and retrieves the values in the scan_acquisition_time and total_intensity variables. Then, with the provided sampling rate and modulation time, the chromatogram is folded into a numerical matrix (two-dimensional chromatogram)	name	The name of the netCDF file which the data will be retrieved	[31]
				mod_time	The modulation time of the chromatographic run	
				sam_rate	The sampling rate of the equipment. If sam_rate is missing, the sampling rate is calculated by the dividing one by the difference of two adjacent scan time	
				per_eval	An integer with the percentage of the run time to be evaluated, if the sampling rate is homogeneous	
				x_cut	The retention time range in the first dimension to be maintained while importing data.	-
				y_cut	The retention time range in the first second to be maintained while importing data.	-
	Plot two-dimensional chromatogram	plot	This function receives a two-dimensional matrix with TIC signals and plots them into a contour or filled contour plot	type	A single character indicating the type of chromatogram representation. By default, type= "f" for a filled contour, if type = "c" only contours or isolines will be displayed.	-
Preprocessing	Smoothing	wsmooth	This function takes a raw two-dimensional chromatogram and performs the weighted Wittaker smoother routine. It smooths with linear or quadratic penalty alongside the first dimension, based on Whittaker smoother	chromatogram	A two-dimensional chromatogram	[33]
				penalty	The penalty order. Only penalty of first (penalty = 1) and second-order (penalty = 2) are allowed. By default, it is performed with the first penalty order.	
				lambda	smoothing parameter: larger values lead to more smoothing	
	Baseline correction	baseline_corr	It corrects the baseline of the chromatogram based on ALS	chromatogram	A two-dimensional chromatogram	[34]

Peak alignment	twod_cow	It aligns a sample chromatograms against a reference chromatogram by two-dimensional correlation optimized warping algorithm	sample_chrom	A two-dimensional chromatogram which will be aligned	[35]
			ref_chrom	A reference two-dimensional chromatogram	
			segments	Two integers with the number of segments in which the first and second dimension will be subdivided, respectively.	
			max_warp	A two integer vector with the maximum warping parameter	
Multiway Principal Component Analysis	m_prcomp	It performs a multiway principal component analysis on the given two-dimensional chromatograms. Before to perform the calculation, each given chromatograms are unfolded to a one-dimensional vector	chrom	Multiple chromatograms or batch ones aligned	[36]
			center	A logical value indicating whether the variables should be shifted to be zero centered. True is set by default and is strongly suggested not to change to False.	
			scale	a logical value indicating whether the variables should be scaled to have unit variance before the analysis takes place. The default is True to give the same variable importance in chemometrics.	
			npcs	An integer indicating how many principals components are desired to maintain. The default is 3 principal component	
Retrieve scores from MPCA	scores	exports the scores matrix of the previously MPCA performed	Object	The result of m_prcomp function	-
Plot MPCA loadings	plot_loading	This function takes the loadings of MPCA and eval if a certain variable was removed previous to compute the MPCA and fill the removed variables with zero (zero variance variables).	Object	The result of m_prcomp function	-
			type	The value type of loadings, “p” for positive, “n” for negative, and “b” for negative and positive loading values	
			pc	The number of the principal component to the plot	

2.2.2. Visualization

As described by Reichenbach [37], the GC×GC image visualization consists of colored pixels layers. In this case, interpolation techniques are used for a two-dimensional image view. The interpolation function evaluates the collection of TIC at the $C_{(I)}$ and $C_{(J)}$ coordinates, and then approximate the contours with the computed interpolation. On the other hand, is common in chromatography that the signals contain highly similar intensities, or a certain group of chemical entities produces signals several times higher than the rest of the molecules. Therefore, in this library, the user has two options to display two-dimensional chromatograms, filled contour, and contour plot. While contour plot displays low and high-intensity TIC signals as isolines with a white background, filled contour assigns a different color to the background, which may obscure low-intensity signals. In other words, analyte concentrations are not usually evenly distributed in the sample, they may produce a large range of signal intensities. Thus, analytes with greater concentrations may opaque the visualization of the analytes with lower concentrations. To overcome this issue, the user can choose the contour plot to display low signal intensities.

Moreover, the color palette must be taken into consideration in this type of graphical representations, since they play the main role by capturing reader attention [38]. Building an effective color ramp may be difficult for inexperienced users. We encourage users to employ the `colorRmp`s package [39]. On the other hand, users with more programming skills can create a color palette from scratch, as explained in [38].

2.2.3. Pre-processing

High throughput chemical equipment achieves a great level of detail, in which external artifacts are also included, such as instrument variability or sample matrix effect. Therefore, the analyst has to eliminate undesirable information to convert the raw data into useful information, because it has a huge effect on the downstream analysis [11,21,29]. Consequently, several pre-processing techniques have been developed in order to remove chemical and instrument noise. In general, pre-processing techniques include three basic modules: smoothing, baseline correction, and peak alignment.

2.2.3.1. Smoothing

Denoising signals enhance the signal to noise ratio (S/N) in the chromatogram, increasing both accuracy and precision analytical results. Although the pioneers in smoothing were Savitzky and Golay with the local likelihood approach [40,41], Whittaker smoother was introduced as a general-purpose algorithm in chemistry, showing better performance [33]. Whittaker smoothing works in the time domain by discrete penalized least squares taking into consideration data fidelity of original data and roughness of the fitted data. This algorithm starts with a supposition of a noisy signal y and a smoothed signal z that fits y . The roughness (R) of the smoother can be described as the sum square of z $R = \sum_i (z_i - z_{i-1})^2$. Moreover, the lack of fitting can

be expressed as the sum of squares of differences $S = \sum_i (y_i - z_i)^2$. Finally, the governing equation of Whittaker smoother can be stated as:

$$Q = S + \lambda R \quad \text{eq. (1)}$$

Where λ is a user given multiplicative factor to the roughness. The aim of Whittaker smoother is to find the combination of z that minimizes Q . While the user gives λ larger values, z will be more smoothed. Different results by varying λ are showed in [42]. These advantages account for automatically boundaries adaptation, missing values and sparse handling, and good computational efficiency in a desktop computer. The Whittaker routing is available through the “wsmooth” function.

2.2.3.2. Baseline correction

The baseline drift in GC/GC×GC is mostly caused by column bleeding or complex mixtures that cannot be separated [11]. In order to remove this type of noise, baseline correction removes the baseline noise and centering the signal around zero. The proposed library implements baseline correction by asymmetric least squares algorithm. Eilers proposed the baseline correction by adapting Whittaker smoothing to calculate the trend of the baseline [43]. In this extension, weights (ω) are introduced, stated as.

$$Q = \sum_i \omega_i (y_i - z_i)^2 + \lambda \sum_i (\Delta^2 z_i)^2 \quad \text{eq. (2)}$$

While ordinary least squares obtain the residual based on the difference of the raw and fitted signal ($y - z$), and the sign of the residuals has the same effects over the penalties, the asymmetric least squares give more weight to negative residuals than positive residuals. This approach is considered based on the positive residuals are obtained when a peak is detected. Therefore, the analytical peak signal has not had to be distorted. In contrast, negative residuals are more penalized. As stated above, weights are assigned based on the sign of the residual as follows: $\omega_i = p$ if $y_i > z_i$ and $\omega_i = 1 - p$. Here, p is introduced as a user parameter. In cases when the fitted signal is greater than the raw signal, weights are the difference of $\omega_i = 1 - p$ as proposed by Eilers [43].

Although asymmetric least squares offer advantages, such as short computational times, parameter flexibility, the parameters have to optimize by hand. This function is available in the library through the “baseline_corr” command.

2.2.3.3. Peak alignment

The retention shift in chromatography is an unavoidable source of experimental variation [21]. Retention shift can be caused by stationary phase decomposition, column change during usage, or different modulation temperatures over the experiment. As described above, this software follows the stream of pixel-based analysis. Therefore, 2DCOW algorithm was implemented [35].

Basically, the 2DCOW works by splitting the sample (A) and reference (B) chromatogram of $X_{(I, J)}$ dimensions into m segments for both dimensions (1D and 2D), respectively. The new partitioned matrix can be stated as n_i and n_j with column nodes in the first $\{e_0, e_1, \dots, e_{n_i}\}$, and the second $\{f_0, f_1, \dots, f_{n_j}\}$ column, which are a segment subset of $\{1, 2, \dots, I\}$ and $\{1, 2, \dots, J\}$, respectively. Each grid from the partitioned matrix, for the sample and the reference chromatogram, can be expressed as $\{(e_k, f_l) : k = 0, 1, \dots, n_i ; l = 0, 1, \dots, n_j\}$. At each row e_k of A, a new row vector $\tilde{A}_{e_k} = (\tilde{A}_{e_k1}, \tilde{A}_{e_k2}, \dots, \tilde{A}_{e_kl})$ is obtained, with the j th component \tilde{A}_{e_kf} through.

$$\tilde{A}_{i_kj} = \sum_{e=1}^n A_{ef} W e^{-e_k/h} / \sum_{e=1}^n W e^{-e_k/h} \quad \text{eq. (3)}$$

This routine requires one pair of arguments for each of the first and second dimensions that are called segment length, which is the number of sections to split the chromatogram and slack, which is the maximum warping level. In the proposed library, there are two functions to perform two-dimensional peak alignment, “twod_cow” and “batch_2DCOW”. While the first command can align a single chromatogram against a reference, the second routine can align a set of chromatograms to a reference. Regarding the reference chromatogram, it has to keep high peak similarity between sample chromatograms since the signals should match to perform the alignment, a more detailed explanation about criteria to choose the target chromatogram is provided in [11,15]. One of the main advantages that provides 2DCOW is the interpolative warping of the warped region and the reference in order to maximize the correlation between them, correcting the retention time shifts.

One strategy to select the reference chromatogram is to create an artificial chromatogram by summing or averaging pixels of a set of chromatograms from different groups. In the proposed toolbox, we provide the “ref_chrom” function, which receives multiple chromatograms and computes a new temporal chromatogram to be employed as the reference.

2.2.4. Multivariate analysis

Multivariate methods are capable of simultaneously analyzing multiple variables in order to expose group-wise variation. There are two flavors for multivariate analysis focused on statistical learning, supervised and unsupervised analysis. Supervised analysis requires prior information about the sample space composition, a predicted variable, which commonly is the sample class, to train the model. In contrast, the unsupervised approach does not need extra information about sample composition to compute its routines and creates a discrimination model. In consequence, supervised approaches receive more information about group arrangement, and usually show better results than unsupervised techniques [44]. Concerning to multivariate analysis, the toolbox presents multiway principal component analysis. Also, RGCxGC presents capabilities to export chromatographic data

in a compatible format to be used in external toolboxes. For example, the exported chromatograms can be subjected to the partial least squares-discriminant analysis available in the mixOmics toolbox [45].

Principal Component Analysis is probably the most unsupervised analysis in many areas with dimensional reduction purposes. It was introduced in 1901 by Karl Pearson [46]. In chemometrics, PCA has been widely applied for pattern recognition. An adaptation of PCA, is MPCA explained by Wold [36], which can deal with higher-order data. Although there are methods that can analyze high-order data, the authors conclude that the same results can be obtained by unfolding the data into two-way matrices [47,48]. In the case of GC×GC, as stated above, each chromatogram consists of a two-way matrix of dimensions $A_{(I, J)}$, being I the number of modulations per run, and J the number of scans per modulation. the unfolding procedure is carried out as follows: the modulation I +1 is concatenated after the modulation I, and so on for all modulations. As a result, the two-dimensional chromatogram has been unfolded into a one-dimensional row-wise vector. All chromatograms are subjected to this procedure in order to obtain a two-way matrix where the columns are the retention times and rows are samples. Then an ordinary PCA can be performed.

The PCA then decomposes $X_{(q,r)}$ into a score (S) and loading (L) matrices, so that $X = SL^T$. Score matrix is the projection in the reduced multivariate space spanned by principal components, and it is related to the (chromatographic) differences among the samples. On another hand, the loading matrix explains the relationship between variables, where positive values refer to similarities between variables and negative values denote differences in variables across samples [49]. The MPCA can be performed with the “m_pcomp” command.

Once data was subjected to MPCA, the user can access the chromatogram projection into the principal component space (scores matrix) by the “scores” command. Differently in the loading matrix, each eigenvector contains all input variables. Then, each principal component should be interpreted as a two-dimensional chromatogram. Thus, each eigenvector is retrieved from MPCA and folded again into a typical GC×GC chromatogram. This task is carried out by the “plot_loading” command. Finally, one extra parameter in MPCA is the explained variance, which can be retrieved through the “print” command.

2.2.4.1. External discriminant analysis.

Although PCA based techniques are one of the most unsupervised algorithms, GC×GC data can also be subjected to supervised algorithms, such as the Partial Least Squares-Discriminant Analysis (PLS-DA) which is one of the most common classification model widely applied in chemometrics [50]. Therefore, in order to ensure a set of discriminant analysis that can be performed once the chromatograms are exported, we also include the PLS-DA analysis. the RGCxGC library can export chromatograms in a friendly structure to communicate with external libraries. In this work, we extend our analysis by testing

our datasets with PLS-DA available in mixOmics [45]. Even though, more external libraries can be used to employ any desired classification algorithm such as hierarchical clustering, artificial neural networks or supporting vector machine.

2.3. Benchmarking

In order to evaluate the performance with large GC×GC-TOF/MS data, we perform a benchmarking of every important algorithm in RGCxGC package with chromatograms from the *Salmonella* dataset, in order to simulate a large scale experiment. The *Salmonella* dataset was chosen since the mass analyzer used in this study was a time of flight, which has higher resolution and also is dimensional higher. We measured the computational time elapsed to perform the desired routine from 1 sample, with an upper limit of 100 samples with increments of 10 samples per step, except in the first increment where it was 9 samples. The benchmarking was performed in a computer with an Intel Core i7 2.7 GHz processor with a Linux based operating system.

2.4. Application

The proposed software was fully tested with a real laboratory experiment based on microbial antagonism and two published datasets. The first datasets came from research focused on the discriminant analysis of chronic typhoid carriers acquired with a TOF mass analyzer [51]. While the second datasets were related to antiphytopathogenic interaction of different yeast strains against *P. digitatum* acquired with quadrupole mass analyzer [52]. The third dataset was previesly described in the method section (2.1.).

One principal advantage to work with the published dataset is the verification of the software performance through the comparison with real datasets, and being able to show the capability to work with different mass analyzers (Tab. 2). This is an important characteristic because a more sophisticated mass analyzer (TOF) provides more resolution, making data analysis the most computational challenging task.

Table 2. Description of the dataset analyzed with the proposed RGCxGC toolbox. While the dimensions of the single chromatogram are given by $A_{(I, J)}$, the dimensions of the entire datasets are given by $A_{(I, J, K)}$ where K represents the number of samples in each experiment.

Datasets	Samples	Number of categories	Chromatogram dimensions $A_{(I, J)}$	Mass Analyzer	Reference
<i>Salmonella</i>	30	3	$A_{(500, 709)}$	TOF	[32]
<i>Penicillium</i>	18	2	$A_{(150, 368)}$	Quadrupole	[52]
<i>Myrothecium</i>	38	5	$A_{(125, 381)}$	Quadrupole	in-house

3. Results and discussion

In the following section, we present the performance of the RGCxGC package with three different datasets collected with different mass analyzers. We have also performed a supervised analysis (PLS-DA) to confirm the connection of RGCxGC with external tools.

3.1. *Salmonella* dataset

The *Salmonella* dataset was downloaded from MetaboLights database with the MTBLS579 identifier [32]. Chromatograms were downloaded and manually checked for consistency; those with different dimensions specified in Tab 2, were removed. The entire data comprises 30 blood plasma samples with two main categories related to carriage diagnosis, control samples, and *Salmonella sp.* carriage. Moreover, in the second category, authors collected samples from two different etiological agents; *S. typhi* and *S. Paratyphi A*. a representative chromatogram, based on the number of peaks detected, is presented in the Figure 2A.

All chromatograms were smoothed with a quadratic penalty and a λ equal to 10. Then, the baseline correction was performed with a correction factor equal to 1000. For the two-dimensional alignment, a sample of confirmed *S. Paratyphi* sample carriage (07_GB, MetaboLight identifier) was chosen to be used as a template and the remaining chromatogram was aligned against it. First and the second dimension was divided into 20 and 40 segments, respectively. The maximum warping values for the first and second dimensions were 2 and 8, respectively. Prior to principal component analysis, chromatograms were mean-centered.

Samples groups do not present a clear classification in the projected principal component space. Also, the explained variance for this dataset is the lowest (< 15%) obtained by MPCA in the first two principal components (Fig. 2B and 2C) and there is not a clear difference between the two etiological agents. This poor differentiation between *S. typhi* and *S. Paratyphi A* carriage patients maybe for the chromatogram similarity and the performance of the statistical learning algorithm. Subsequently, chromatograms were subjected to a supervised learning algorithm, PLS-DA. For this analysis, the aligned chromatograms were exported by unfolding them into a matrix. This matrix represents the explanatory variables, while the carriage status represents the predicted variable. The PLS-DA was performed according to the mixOmics procedure [45]. In contrast with the MPCA results, the PLS-DA model defines clusters between control and carriage patients (Fig. 2D). Furthermore, the model also clusters the two different etiological agents (*S. paratyphi A* and *S. typhi*). The classification improvement can be explained for the type of statistical learning employed. In the case of PLS-DA, it is a supervised algorithm, which is trained with the correct sample category. It is not surprising since supervised algorithms show better performance than unsupervised algorithms. Even though, the explained variance is not higher than the 15%.

3.2. *Penicillium* dataset

Raw chromatograms were provided by the authors of the anti-phytopathogenic yeast strains experiment [53]. *Penicillium digitatum* is a pathogen that infects citrus fruits, causing product degradation within the product's storage, transportation, and market. For this reason, microbial activity was tested against different yeast strains by a coculture experiment in triplicates. Yeast strains, with validated high and low antagonist activity, were selected for downstream analysis.

All chromatograms were baseline corrected with a correction factor equal to 0.5. Then, chromatograms were smoothed with a linear penalty and a λ equal to 2 (Fig. 2E). For the peak alignment process, a consensus chromatogram was computed by averaging the pixel values for multiple chromatograms, as explained above, and the remaining chromatogram was aligned against it. The maximum warping value for the first and second dimension were 4 and 10, respectively. Prior to principal component analysis, chromatograms were mean-centered.

The separation achieved by MPCA was clearly notorious. In this context, the total explained variance between the first two PC was 60.48% (Fig. 2F and 2G). While all low activity yeast strains are clustering at the negative PC1 scores, all high active antagonist yeast strains clustered in positive values. Furthermore, triplicates of each high active yeast strains were clustered together, expressing replicate and strain similarities. For example, the strain pe2 (*Saccharomyces cerevisiae* ACB-PE2) has different profiles than cat1 and kd1, (*S. cerevisiae* ACB-CAT1 and *S. cerevisiae* ACB-KD1) together.

Furthermore, a PLS-DA classification was also conducted, as a supervised approach. For this analysis, the aligned chromatograms were exported by unfolding them into a matrix. This matrix represents the explanatory variables, with anti-phytopathogenic activity being the predicted variable. In this case, the MPCA and PLS-DA results show high similarities, both models can differentiate high and low strain activities in the first projected dimension. The PLS-DA was not able to separate yeast strains with low antagonist activity, it classifies all strains into a single cluster (Fig. 2H).

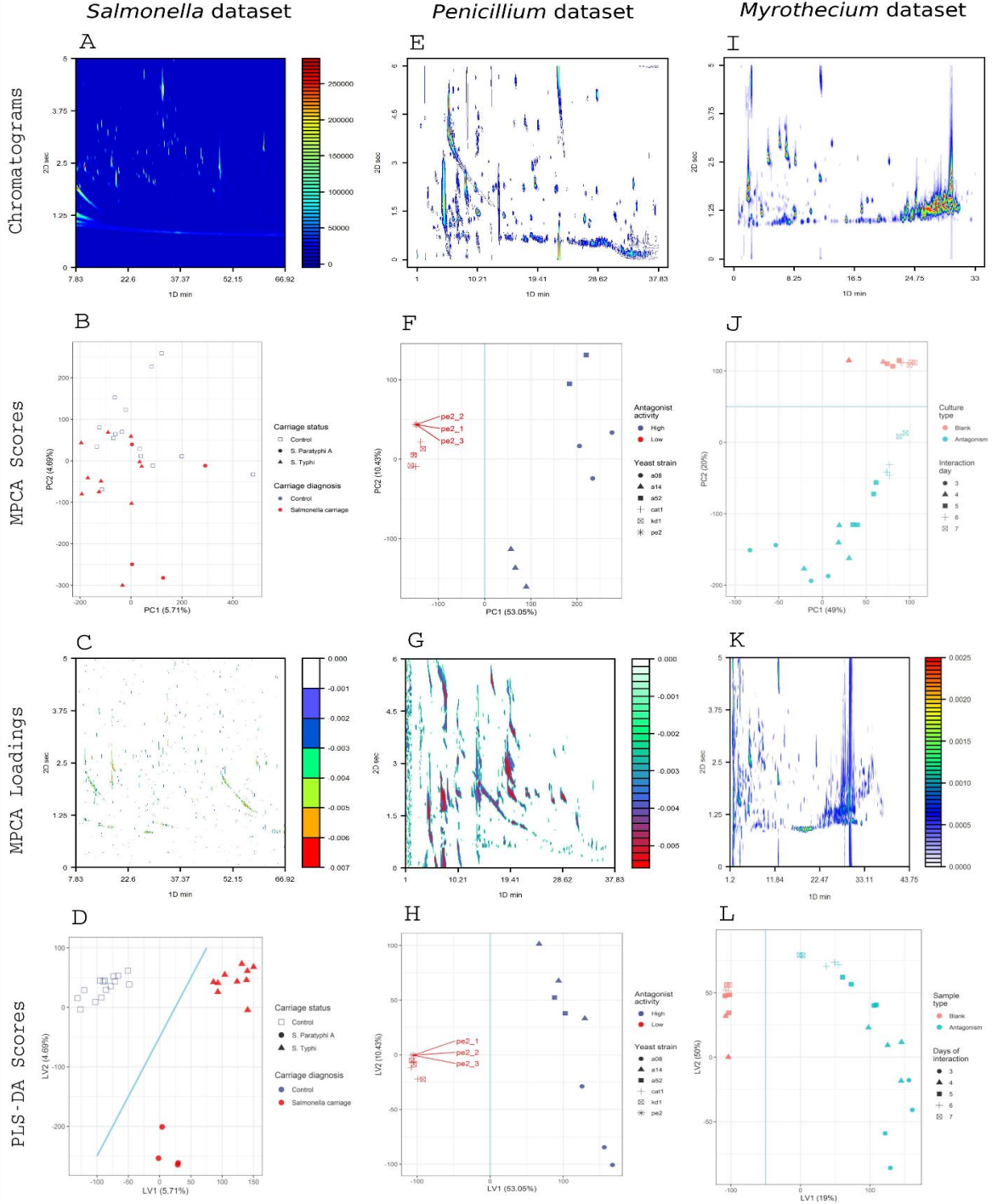


Figure 2. The grid shows the results of the three datasets analyzed with the proposed toolbox; *Salmonella* (A, B, C, and D), *Penicillium* (E, F, G, and H) and *Myrothecium* (I, J, K, L). The first row shows the representative two-dimensional TIC chromatograms, the second row presents the MPCA scores, the third row displays the MPCA loadings, and the fourth row displays the PLS-DA scores. The results of unsupervised discriminant analysis by MPCA, of the *Penicillium* (F) and *Myrothecium* (G) datasets, present a clear separation between categories (edges in blue) in the two first PC's, with an explained variance greater than 50%. On the other hand, the MPCA cannot explain the differences in the *Salmonella* (B) dataset categories. In the case of supervised classification with PLS-DA, the model can discriminate between categories for all datasets (D, H, and L).

3.3. *Myrothecium* dataset

The *Myrothecium* dataset was locally made with the aforementioned methodology in the method section (2.2). The *Myrothecium* genus plays a major role for combating white mold in soybeans. Therefore, the MYL volatilome kinetics were explored. The entire data comprise 38 samples with two main categories, control (Bco) and *Myrothecium* (MYL) culture kinetics. Within the MYL category, there were five subcategories in concordance with days that the fungus has grown in the culture media.

For the peak alignment process, the chromatogram from the MYL culture on the fifth day was selected as a representative chromatogram, and the remaining chromatograms were aligned against it (Fig. 2I). The maximum warping value for the first and second dimension were 4 and 10, respectively. Prior MPCA, chromatograms were mean-centered.

In the MPCA score plot, the separation between control and antagonism samples was clearly appreciated (Fig. 2J and 2K). Thus, the accumulated explained variance in the two first principal components of this model was about 70%. In contrast with the two previously discussed datasets, in this case, groups were separated by the second PC. Therefore, for metabolite annotation, compounds with the highest eigenvector values were annotated, resulting in 48 metabolites with higher values of similarity than 80% in the database (see Annexed). Between the annotated metabolites, several chemical species were found that were previously described to have fungal biocontrol activity. For example, the presence of phenylethyl alcohol, α -curcumenol and α -terpinene that were described in *Trichoderma* genus control [54]. In the same manner, these compounds had also been reported to be produced by *Memnoniella* genus, which can induce pathogen resistance in plants through its volatilome [55]. Moreover, phenylethyl alcohol was also found in antimicrobial extracts from the *Glicladium* genus [56].

Finally, we also performed a PLS-DA model on the in-house *Myrothecium* dataset. As stated earlier, the pre-processed chromatograms were exported by unfolding them into a two-way matrix. The PLS-DA, also present consistency with the MPCA, with a rotation in the latent variables (Fig. 2L). In other words, in contrast with the MPCA that captures the variance in the second projected variable (PC2), the PLS-DA captured the experiment variance in the first projected variable (latent variable). Moreover, both models can also explain intra-day antagonism variability.

3.4. Computational efficiency

In order to simulate a large scale experiment, we performed the benchmarking of the main functions of the proposed toolbox. Then, chromatograms of the *Salmonella* dataset were chosen due to the higher dimensions, as described in Tab 1. All available functions showed to have a linear increment in the elapsed time to be performed (Fig. 3). Moreover, the longest time required to perform a chemometric routine with 100 samples was 82 minutes was the “twod_cow” routine. This was expected since peak

alignment is the most time-consuming task and the major bottleneck in this type of analysis, since different algorithms and parameters may be tested before useful information can be extracted from the raw signals. In the case of the MPCA routine, the benchmarking showed to have the second longest elapsed time, requiring 17.5 min to analyze 100 samples. Furthermore, functions that manage graphical components (“plot” and “plot_loading”) require over 3 minutes to display 100 samples. On the other hand, the rest of the functions did not need half of a minute to work with the maximum number of samples tested. In comparison with similar alignment methodologies for GC×GC, this procedure is the most time-consuming routine [57]. For example, in the case of a MATLAB pixel-by-pixel alignment, between 10 to 20 min were required to align a 1 sample [57]. In the case of RGCxGC package, with the same time range, 10 to 20 samples can be processed.

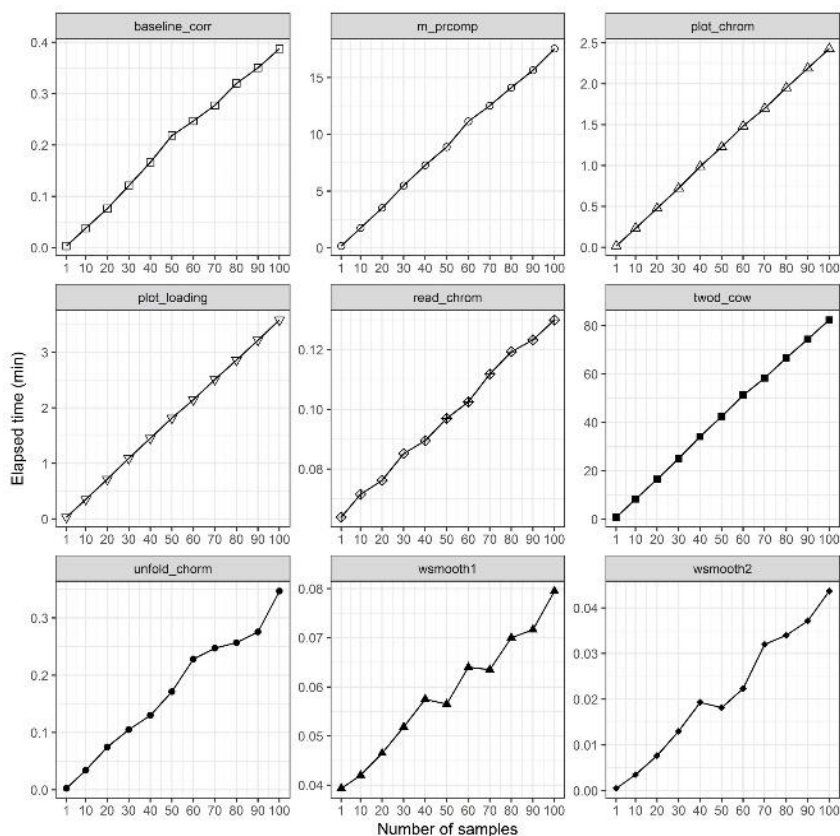


Figure 3. Benchmarking results of the main functionalities of the RGCxGC toolbox. The head of each figure represents the function name that was tested. In order to create a similar situation of a big scale experiment, we tested the elapsed time for each function from 1 to 100 chromatograms samples in increments of 10 samples, except for the first step where the increment was 9 samples. All functions present a linear increment while the number of samples increases. Meanwhile “m_prcom” and “twod_cow” require 20 and 82 minutes to process 100 samples, the rest of the functions requires less than 5 minutes to process 100 samples.

4. Conclusions

In this manuscript, we presented a novel end-to-end workflow for non-targeted GC×GC-MS exploratory data analysis by signal processing and with statistical learning algorithms. The currently available functions for signal processing in the RGCxGC package are compiled with baseline correction, smoothing, two-dimensional peak alignment. While for statistical learning, the multiway principal component analysis was implemented for unsupervised classification and partial least squares discriminant analysis was tested. Furthermore, we provide a generic manner to connect with other libraries in order to provide a wide spectrum of possible classification algorithms, such as linear discriminant analysis, artificial neuronal networks. The presented software showed to be capable to process a considerable amount of data (> 1 GB). Also, the longest required time for the desired routine to be performed was less than 2 hours. This is an advantage for large-scale experiments, such as metabolomics studies, to overcome the bottleneck of data analysis. On the other hand, the characteristic of free open source implementation could help with research reproducibility and reduce the dependence of private license depending software. Our approach was successfully tested in two published datasets and one in-house dataset. The key benefit of this implementation is to avoid multiple software in non-target studies. In addition, the software is continuously checked and maintained by package developers and CRAN team, in order to ensure long term user availability and avoid obsolescence. The proposed library, as well as a detailed user manual and a complete tutorial, is freely available and can be found at <https://cran.r-project.org/web/packages/RGCxGC/index.html>.

References

- [1] C. Poole, *Gas Chromatography*, 1st ed., 2012.
- [2] S.E. Prebihalo, K.L. Berrier, C.E. Freye, H.D. Bahaghighat, N.R. Moore, D.K. Pinkerton, R.E. Synovec, *Multidimensional Gas Chromatography: Advances in Instrumentation, Chemometrics, and Applications*, *Anal. Chem.* 90 (2018) 505–532. <https://doi.org/10.1021/acs.analchem.7b04226>.
- [3] P.Q. Tranchida, *Comprehensive two-dimensional gas chromatography: A perspective on processes of modulation*, *J. Chromatogr. A.* 1536 (2018) 2–5. <https://doi.org/10.1016/j.chroma.2017.04.039>.
- [4] B. Gruber, B.A. Weggler, R. Jaramillo, K.A. Murrell, P.K. Piotrowski, F.L. Dorman, *Comprehensive two-dimensional gas chromatography in forensic science: A critical review of recent trends*, *TrAC - Trends Anal. Chem.* 105 (2018) 292–301. <https://doi.org/10.1016/j.trac.2018.05.017>.
- [5] A.M. Muscalu, T. Górecki, *Comprehensive two-dimensional gas chromatography in environmental analysis*, *TrAC - Trends Anal. Chem.* 106 (2018) 225–245. <https://doi.org/10.1016/j.trac.2018.07.001>.
- [6] B.J. Pollo, G.L. Alexandrino, F. Augusto, L.W. Hantao, *The impact of comprehensive two-dimensional gas chromatography on oil & gas analysis: Recent advances and applications in petroleum industry*, *TrAC - Trends Anal. Chem.* 105 (2018) 202–217. <https://doi.org/10.1016/j.trac.2018.05.007>.
- [7] T. Miyazaki, K. Okada, T. Yamashita, M. Miyazaki, *Two-dimensional gas chromatography time-of-flight mass spectrometry-based serum metabolic fingerprints of neonatal calves before and after first colostrum ingestion*, *J. Dairy Sci.* 100 (2017) 4354–4364. <https://doi.org/10.3168/jds.2017-12557>.
- [8] J.C. Giddings, *Sample dimensionality: A predictor of order-disorder in component peak distribution in multidimensional separation*, *J. Chromatogr. A.* 703 (1995) 3–15. [https://doi.org/10.1016/0021-9673\(95\)00249-M](https://doi.org/10.1016/0021-9673(95)00249-M).
- [9] R.B. Wilson, W.C. Siegler, J.C. Hoggard, B.D. Fitz, J.S. Nadeau, R.E. Synovec, *Achieving high peak capacity production for gas chromatography and comprehensive two-dimensional gas chromatography by minimizing off-column peak broadening*, *J. Chromatogr. A.* 1218 (2011) 3130–3139. <https://doi.org/10.1016/j.chroma.2010.12.108>.
- [10] L. Yi, N. Dong, Y. Yun, B. Deng, D. Ren, S. Liu, Y. Liang, *Chemometric methods in data processing of mass spectrometry-based metabolomics: A review*, *Anal. Chim. Acta.* 914 (2016) 17–34. <https://doi.org/10.1016/j.aca.2016.02.001>.
- [11] K.M. Pierce, B. Kehimkar, L.C. Marney, J.C. Hoggard, R.E. Synovec, *Review of chemometric analysis techniques for comprehensive two dimensional separations data*, *J. Chromatogr. A.* 1255 (2012) 3–11. <https://doi.org/10.1016/j.chroma.2012.05.050>.
- [12] E. Szymańska, *Modern data science for analytical chemical data – A comprehensive review*, *Anal. Chim. Acta.* 1028 (2018). <https://doi.org/10.1016/j.aca.2018.05.038>.
- [13] R. Brereton, *Chemometrics for Pattern Recognition*, Wiley, 2009.
- [14] J. Krumsiek, F.J. Theis, *Statistical methods for the analysis of high-throughput metabolomics data Abstract : Metabolomics is a relatively new high-throughput technology that aims at measuring all endogenous metabolites within a biological sample in an unbiased fashion . The resu*, *Comput. Struct. Biotechnol. J.* 4 (2013).
- [15] L. Komsta, Y. Vander Heyde, J. Sherman, *Chemometrics in Chromatography*, 1st ed., CRC Press, 2018.
- [16] D. Zhang, X. Huang, F.E. Regnier, M. Zhang, *Two-dimensional correlation optimized warping algorithm for aligning GCxGC-MS data*, *Anal. Chem.* 80 (2008) 2664–2671. <https://doi.org/10.1021/ac7024317>.
- [17] P.H.C. Eilers, *Parametric Time Warping*, *Anal. Chem.* 76 (2004) 404–411. <https://doi.org/10.1021/ac034800e>.
- [18] B. Wang, A. Fang, J. Heim, B. Bogdanov, S. Pugh, M. Libardoni, X. Zhang, *DISCO: Distance and Spectrum Correlation Optimization Alignment for Two Dimensional Gas Chromatography Time-of- Flight Mass Spectrometry-based Metabolomics*, *Anal. Chem.* 82 (2011) 5069–5081. <https://doi.org/10.1021/ac100064b>.DISCO.
- [19] B. Wang, A. Fang, X. Shi, S.H. Kim, X. Zhang, *DISCO2 : A Comprehensive Peak Alignment Algorithm for Two-Dimensional Gas Chromatography Time-of-Flight Mass Spectrometry*, (2012) 486–491.
- [20] N. Hoffmann, M. Wilhelm, A. Doebbe, K. Niehaus, J. Stoye, *BiPACE 2D-Graph-based multiple alignment for comprehensive 2D gas chromatography-mass spectrometry*, *Bioinformatics.* 30 (2014) 988–995. <https://doi.org/10.1093/bioinformatics/btt738>.
- [21] L. Yi, N. Dong, Y. Yun, B. Deng, D. Ren, S. Liu, Y. Liang, *Chemometric methods in data processing of mass spectrometry-based metabolomics: A review*, *Anal. Chim. Acta.* 914 (2016) 17–34. <https://doi.org/10.1016/j.aca.2016.02.001>.
- [22] S. Castillo, I. Mattila, J. Miettinen, M. Orešič, T. Hyötyläinen, *Data analysis tool for comprehensive two-dimensional gas chromatography/time-of-flight mass spectrometry*, *Anal. Chem.* 83 (2011) 3058–3067. <https://doi.org/10.1021/ac103308x>.
- [23] E.A. Higgins Keppler, C.L. Jenkins, T.J. Davis, H.D. Bean, *Advances in the application of comprehensive two-*

- dimensional gas chromatography in metabolomics, *TrAC Trends Anal. Chem.* 109 (2018) 275–286. <https://doi.org/10.1016/J.TRAC.2018.10.015>.
- [24] R.C. Ramaker, E. Gordon, S.J. Cooper, R2DGC : Threshold-free peak alignment and identification for 2D gas chromatography mass spectrometry in R, *Bioinformatics.* 34 (2017) 1789–1791.
- [25] S. Moayedpour, H. Parastar, RMet: An automated R based software for analyzing GC-MS and GC×GC-MS untargeted metabolomic data, *Chemom. Intell. Lab. Syst.* 194 (2019) 103866. <https://doi.org/10.1016/j.chemolab.2019.103866>.
- [26] R.G. Brereton, G.R. Lloyd, Partial least squares discriminant analysis: Taking the magic away, *J. Chemom.* 28 (2014) 213–225. <https://doi.org/10.1002/cem.2609>.
- [27] Y. Zushi, S. Hashimoto, Direct Classification of GC × GC-Analyzed Complex Mixtures Using Non-Negative Matrix Factorization-Based Feature Extraction, *Anal. Chem.* 90 (2018) 3819–3825. <https://doi.org/10.1021/acs.analchem.7b04313>.
- [28] P.E. Sudol, D. V. Gough, S.E. Prebhalo, R.E. Synovec, Impact of data bin size on the classification of diesel fuels using comprehensive two-dimensional gas chromatography with principal component analysis, *Talanta.* 206 (2020) 120239. <https://doi.org/10.1016/J.TALANTA.2019.120239>.
- [29] Y. Izadmanesh, E. Garreta-Lara, J.B. Ghasemi, S. Lacorte, V. Matamoros, R. Tauler, Chemometric analysis of comprehensive two dimensional gas chromatography–mass spectrometry metabolomics data, *J. Chromatogr. A.* 1488 (2017) 113–125. <https://doi.org/10.1016/j.chroma.2017.01.052>.
- [30] V. Den Dool, P.D. Kratz, A generalization of the retention index system including linear temperature programmed gas-liquid, *J. Chromatography A.* (1962) 463–471. https://doi.org/10.1007/978-3-319-70262-9_7.
- [31] T. Skov, R. Bro, Solving fundamental problems in chromatographic analysis, *Anal. Bioanal. Chem.* 390 (2008) 281–285. <https://doi.org/10.1007/s00216-007-1618-z>.
- [32] E. Näsström, P. Jonsson, A. Johansson, S. Dongol, A. Karkey, B. Basnyat, N. Tran Vu Thieu, T. Trinh Van, G.E. Thwaites, H. Antti, S. Baker, Diagnostic metabolite biomarkers of chronic typhoid carriage, *PLoS Negl. Trop. Dis.* 12 (2018) 1–15. <https://doi.org/10.1371/journal.pntd.0006215>.
- [33] P.H.C. Eilers, A perfect smoother, *Anal. Chem.* 75 (2003) 3631–3636. <https://doi.org/10.1021/ac034173t>.
- [34] J. Peng, S. Peng, A. Jiang, J. Wei, C. Li, J. Tan, Asymmetric least squares for multiple spectra baseline correction, *Anal. Chim. Acta.* 683 (2010) 63–68. <https://doi.org/10.1016/j.aca.2010.08.033>.
- [35] Dabao Zhang, Xiaodong Huang, Fred E. Regnier, Min Zhang, Two-dimensional correlation optimized warping algorithm for aligning GC×GC-MS data, *Anal. Chem.* 80 (2008) 2664–2671. <https://doi.org/10.1021/ac7024317>.
- [36] S. Wold, P. Geladi, K. Esbensen, J. Öhman, Multi-way principal components- and PLS-analysis, *J. Chemom.* 1 (1987) 41–56. <https://doi.org/10.1002/cem.1180010107>.
- [37] S.E. Reichenbach, M. Ni, V. Kottapalli, A. Visvanathan, Information technologies for comprehensive two-dimensional gas chromatography, *Chemom. Intell. Lab. Syst.* 71 (2004) 107–120. <https://doi.org/10.1016/j.chemolab.2003.12.009>.
- [38] E. Pante, B. Simon-Bouhet, marmap: A Package for Importing, Plotting and Analyzing Bathymetric and Topographic Data in R, *PLoS One.* 8 (2013) 6–9. <https://doi.org/10.1371/journal.pone.0073051>.
- [39] K. T. colorRapms, R Packag. (2012).
- [40] T. Hastie, R. Tibshirani, Generalized additive models, *Monogr. Stat. Appl. Probab.* (1990) 15.
- [41] F. J. G. I, Local Polynomial Modelling and Its Applications, CHAPMAN AND HALL, n.d.
- [42] S. Brown, L. Sarabia, T. Johan, Comprehensive chemometrics: chemical and biochemical data analysis, Elsevier, 2009.
- [43] P.H.C. Eilers, Parametric Time Warping, *Anal. Chem.* 76 (2004) 404–411.
- [44] G. James, D. Witten, T. Hastie, R. Tibshirani, An introduction to statistical learning : with applications in R, 2013.
- [45] F. Rohart, B. Gautier, A. Singh, K.-A. Lê Cao, mixOmics: An R package for 'omics feature selection and multiple data integration., *PLoS Comput. Biol.* 13 (2017) e1005752. <https://doi.org/10.1371/journal.pcbi.1005752>.
- [46] K. Pearson, On lines and planes of closest fit to systems of points in space, *London, Edinburgh, Dublin Philos. Mag. J. Sci.* 2 (1901) 559–572. <https://doi.org/10.1080/14786440109462720>.
- [47] L.S. Ramos, K.R. Beebe, W.P. Carey, S.M. Eugenio, B.C. Erickson, B.E. Wilson, B.R. Kowalski, L.E. Wangen, Chemometrics, *Anal. Chem.* 58 (1986) 294–315. <https://doi.org/10.1021/ac00296a020>.
- [48] P. Geladi, Analysis of multi-way (multi-mode) data, *Chemom. Intell. Lab. Syst.* 7 (1989) 11–30. [https://doi.org/10.1016/0169-7439\(89\)80108-X](https://doi.org/10.1016/0169-7439(89)80108-X).
- [49] M. Castro, Quimiometria: Conceitos, Métodos e Aplicações, 1st ed., UNICAMP, Brazil, 2009.
- [50] J. Bartel, J. Krumsiek, F.J. Theis, Statistical Methods for the Analysis of High-Throughput Metabolomics Data,

Comput. Struct. Biotechnol. J. 4 (2013) e201301009. <https://doi.org/10.5936/csbj.201301009>.

- [51] E. Näsström, P. Jonsson, A. Johansson, S. Dongol, A. Karkey, B. Basnyat, N. Tran Vu Thieu, T. Trinh Van, G.E. Thwaites, H. Antti, S. Baker, Diagnostic metabolite biomarkers of chronic typhoid carriage, *PLoS Negl. Trop. Dis.* 12 (2018) 1–15. <https://doi.org/10.1371/journal.pntd.0006215>.
- [52] J.R.B. de Souza, K.C. Kupper, F. Augusto, In vivo investigation of the volatile metabolome of antiphytopathogenic yeast strains active against *Penicillium digitatum* using comprehensive two-dimensional gas chromatography and multivariate data analysis, *Microchem. J.* 141 (2018) 204–209. <https://doi.org/10.1016/j.microc.2018.05.036>.
- [53] J.R. Belinato, K.C. Kupper, F. Augusto, In vivo investigation of the volatile metabolome of antiphytopathogenic yeast strains active against *Penicillium digitatum* using comprehensive two-dimensional gas chromatography and multivariate data analysis, *Microchem. J.* 141 (2018) 362–368. <https://doi.org/10.1016/J.MICROC.2018.05.047>.
- [54] N. Stoppacher, B. Kluger, S. Zeilinger, R. Krska, R. Schuhmacher, Identification and profiling of volatile metabolites of the biocontrol fungus *Trichoderma atroviride* by HS-SPME-GC-MS, 81 (2010) 187–193. <https://doi.org/10.1016/j.mimet.2010.03.011>.
- [55] P.F. De Lima, M.F. Furlan, F.A. De Lima Ribeiro, S.F. Pascholati, F. Augusto, In vivo determination of the volatile metabolites of saprotroph fungi by comprehensive two-dimensional gas chromatography, *J. Sep. Sci.* 38 (2015) 1924–1932. <https://doi.org/10.1002/jssc.201401404>.
- [56] K. Liouane, D. Saïdana, H. Edziri, S. Ammar, J. Chriaa, M.A. Mahjoub, K. Said, Z. Mighri, Chemical composition and antimicrobial activity of extracts from *Gliocladium* sp. growing wild in Tunisia Kaouthar Liouane, *Med. Chem. Res.* 19 (2010) 743–756. <https://doi.org/10.1007/s00044-009-9227-3>.
- [57] Y. Zushi, J. Gros, Q. Tao, S.E. Reichenbach, S. Hashimoto, J.S. Arey, Pixel-by-pixel correction of retention time shifts in chromatograms from comprehensive two-dimensional gas chromatography coupled to high resolution time-of-flight mass spectrometry, *J. Chromatogr. A.* 1508 (2017) 121–129. <https://doi.org/10.1016/j.chroma.2017.05.065>.

Appendix

Supplemental table. Tentative identification of detected metabolites by *Myrothecium sp.* with spectral similarity > 80%.

#	¹ t _r (min)	² t _r (min)	Match (%)	LTPRI _{Lit}	LTPRI _{Exp}	Compound	CAS
1	6.533	0.007	92	926	943	3,7-dimetil-1,6-octadieno (R-citroneleno)	-
2	7.95	0.025	84	993	994	4-ciclohepten-1-ona	-
3	8.367	0.016	91	1008	1007	Acetado de 3-Hexenila	3681-71-8
4	8.867	0.009	94	1018	1020	4-metil-1-(1-metiletil)-1,3-ciclohexadieno (α -terpineno)	99-86-5
5	9.117	0.008	88	1023	1027	1-metil-4-(1-metiletil)-ciclohexeno (para-ment-1-eno)	5502-88-5
6	9.117	0.012	94	1025	1027	1-metil-4-(1-metiletil)-benzeno (para-cimeno)	99-87-6
7	9.2	0.033	97	1030	1040	2-etil-1-hexanol	104-76-7
8	9.283	0.01	85	1046	1042	3,4-dimetil-1,5-ciclooctadieno	21284-05-9
9	11.367	0.011	83	1117	1109	(3E)-4-etil-3-nonen-5-ino	74685-67-9
10	11.867	0.031	81	1125	1122	3,7-dimetil-6-octenal (β -citronelal)	106-23-0
11	12.617	0.068	97	1136	1142	feniletil álcool	60-12-8
12	13.2	0.033	89	1158	1155	(2Z)-6-metil-2-undeceno	-
13	14.033	0.04	81	1181	1176	4-etil-benzaldeído	4748-78-1
14	17.45	0.033	92	1258	1249	1-decanol	112-30-1
15	18.783	0.015	86	1280	1276	1-(eteniloxi)-decano	765-05-9
16	19.45	0.015	84	1299	1289	5-butil-4-noneno	7367-38-6
17	21.617	0.031	91	1351	1341	1,1-dimetil-2-nonilciclopropano	41977-38-2
18	22.2	0.007	89	1357	1355	1-undecanol	112-42-5
19	22.45	0.008	82	1365	1361	undec-(8Z)-enal	-
20	22.783	0.007	83	1376	1369	2-dodecanol	10203-28-8
21	23.45	0.007	89	1384	1386	4,8-dimetiltridecano	55030-62-1
22	23.617	0.009	84	1393	1390	2-butil-1-octanol	2/8/3913
23	23.783	0.011	90	1403	1394	diepi- α -cedreno	469-61-4
24	24.117	0.013	95	1398	1403	[1S-(1 α ,2 β ,4 β)]-2,4-diisopropenil-1-metil-1-vinilciclohexano (β -elemeno)	-
25	24.2	0.008	93	1409	1405	1,1,3-trimetil-2-(3-metilpentil)ciclohexano	-
26	24.617	0.009	86	1411	1415	1-tetradecino	765-10-6
27	24.867	0.008	91	1421	1421	(E)-3-tetradeceno	41446-68-8
28	25.117	0.011	82	1418	1428	α -santaleno	-
29	25.2	0.012	82	1424	1430	β -barbateno	-
30	25.617	0.012	84	1433	1440	α -bulneseno	3691-11-0

31	26.033	0.012	83	1458	1451	(3Z,6E)-3,7,11-trimetil-1,3,6,10-dodecatetraeno (α -farneseno)	26560-14-5
32	26.117	0.009	84	1443	1453	propanoato de citronelil	141-14-0
33	26.533	0.009	85	1466	1463	9a-metildecahidro-2H-benzo[a]ciclohepten-2-ona	-
34	26.7	0.012	81	1475	1467	4-metileno-1-metil-2-(2-metil-1-propen-1-il)-1-vinil- cicloheptano	-
35	26.95	0.012	85	1472	1474	Biciclo[10.1.0]tridec-1-eno	-
36	27.117	0.013	82	1485	1478	8,8-dimetil-9-metileno-1,5-cicoundecadieno	-
37	27.283	0.016	86	1480	1482	1-(1,5-dimetil-4-hexenil)-4-metil-benzeno (α - curcumeno)	644-30-4
38	27.7	0.01	84	1483	1492	isobutirato de citronelil	97-89-2
39	27.867	0.012	84	1503	1496	1-metil-4-(4-metilpentil)-3-ciclohexeno-1- carboxaldeido	-
40	27.95	0.013	86	1499	1499	2,6,10,10-tetrametilbiciclo[7.2.0]undeca-2,6-dieno (cariofileino-(I3))	-
41	28.117	0.011	84	1501	1503	butirato de citronelil	141-16-2
42	28.7	0.013	87	1518	1517	cis- α -bisaboleno	29837-07-8
43	30.033	0.012	81	1557	1551	(1E,5E)-1,5-dimetil-8-(propan-2- ilideno)ciclodeca-1,5-dieno	15423-57-1
44	30.367	0.013	83	1567	1559	6,10-dimetil-dodeca-1,6-dien-12-ol	-
45	30.7	0.011	82	1574	1567	11-tridecin-1-ol	33925-75-6
46	32.95	0.027	80	1635	1632	(4,8,8-trimetildecahidro-1,4-metanoazulen-9- il)metanol (Isolongifolol)	1139-17-9
47	34.283	0.025	81	1683	1676	1,3-di(ciclohexil)but-1-eno	-
48	35.117	0.028	87	1706	1704	4-isopropil-1,6-dimetilnaftaleno (cadaleno)	483-78-3