

MALDI imaging mass spectrometry and chemometric tools to discriminate highly similar colorectal cancer tissues

S. Mas, A. Torro, L. Fernández, N. Bec, C. Gongora, C. Larroque, P. Martineau, A. de Juan, S. Marco

▶ To cite this version:

S. Mas, A. Torro, L. Fernández, N. Bec, C. Gongora, et al.. MALDI imaging mass spectrometry and chemometric tools to discriminate highly similar colorectal cancer tissues. Talanta, 2020, 208, pp.120455. 10.1016/j.talanta.2019.120455 . hal-03003933

HAL Id: hal-03003933 https://hal.science/hal-03003933

Submitted on 28 Dec 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

MALDI imaging mass spectrometry and chemometric tools to discriminate highly similar colorectal cancer tissues.

S. Mas^{1,2*}, A. Torro³, L. Fernández^{1,4}, N. Bec^{3,5}, C. Gongora³, C. Larroque^{3,6}, P. Martineau^{3¥}, A. de Juan^{2¥} and S. Marco^{1,4¥}

¹Signal and Information Processing for Sensing Systems, Institute for Bioengineering of Catalonia (IBEC), The Barcelona Institute of Science and Technology, Baldiri Reixac 10-12, 08028 Barcelona, Spain

²Chemometrics Group. Department of Chemical Engineering and Analytical Chemistry. UB. Av. Diagonal, 645. 08028 Barcelona, Catalonia, Spain

³Institut de Recherche en Cancérologie de Montpellier (IRCM), INSERM U1194, Université de Montpellier, Institut régional du Cancer de Montpellier (ICM), Montpellier, F-34298, France.

⁴Institute for Regenerative Medicine & Biotherapy (IRMB), INSERM U1183, CHRU of Montpellier, 80 Rue Augustin Fiche, Montpellier, F-34295, France

⁵Department of Electronics and Biomedical Engineering, Universitat de Barcelona,

Marti i Franqués 1, Barcelona 08028, Spain.

⁶Supportive Care Unit, Institut du Cancer de Montpellier (ICM), 208 Rue des Apothicaires, Montpellier, F-34298, France

¥ These authors share senior authorship.

*Corresponding author: Sílvia Mas Garcia, Signal and Information Processing for Sensing Systems, Institute for Bioengineering of Catalonia (IBEC), Baldiri Reixac 4-8, Barcelona 08028, Spain. Telf: +34 934 031118. <u>silviamasgarc@hotmail.com</u>

Abstract

Intratumour heterogeneity due to cancer cell clonal evolution and microenvironment composition and tumor differences due to genetic variations between patients suffering of the same cancer pathology play a crucial role in patient response to therapies. This study is oriented to show that matrix-assisted laser-desorption ionization-Mass spectrometry imaging (MALDI-MSI), combined with an advanced multivariate data processing pipeline can be used to discriminate subtle variations between highly similar colorectal tumors.

To this aim, experimental tumors reproducing the emergence of drug-resistant clones were generated in athymic mice using subcutaneous injection of different mixes of two isogenic cell lines, the irinotecan-resistant HCT116-SN50 (R) and its sibling human colon adenocarcinoma sensitive cell line HCT116 (S). Because irinotecan-resistant and irinotecan-sensitive are derived from the same original parental HCT116 cell line, their genetic characteristics and molecular compositions are closely related.

The multivariate data processing pipeline proposed relies on three steps: (a) multiset multivariate curve resolution (MCR) to separate biological contributions from background; (b) multiset K-means segmentation using MCR scores of the biological contributions to separate between tumor and necrotic parts of the tissues; and (c) partial-least squares discriminant analysis (PLS-DA) applied to tumor pixel spectra to discriminate between R and S tumor populations. High levels of correct classification rates (0.85), sensitivity (0.92) and specificity (0.77) for the PLS-DA classification model were obtained. If previously labeled tissue is available, the multistep modeling strategy proposed constitutes a good approach for the identification and characterization of highly similar phenotypic tumor subpopulations that could be potentially applicable to any kind of cancer tissue that exhibits substantial heterogeneity.

Keywords

MALDI imaging, tumor heterogeneity, chemometrics, multivariate analysis, colorectal cancer

Introduction

Colorectal cancer (CRC) is the third cause of cancer mortality worldwide, affecting men and women almost equally. Despite significant improvements in early cancer detection and treatment, about half of the patient will develop distant metastases (mCRC) associated with a 5-year survival rate for patients of less than 10 % [1]. The standard of care for mCRC includes surgery of the primary tumor, if still present, combined with chemotherapy based on Oxaliplatin, Irinotecan and 5-Fluorouracile administered in various combinations and an antibody-based chemotherapy. These treatments may induce regression of metastases and, in favorable cases, a secondary surgery of liver metastases.

Tumor heterogeneity is now recognized as a key driver in patient response to therapy and thus in clinical outcome. Whereas intertumor heterogeneity among patients has been identified for many years as an important factor explaining differences in patient response to therapeutic regimens, intratumor heterogeneity and its implication in resistance mechanisms has been only more recently recognized and addressed [2–7]. Understanding and characterizing tumor heterogeneity is therefore a key factor to improve patient management and treatment.

To characterize these heterogeneities, a technique that maintains the spatial organization of cells but allows in depth analysis of their molecular content in eventually unknown molecules is needed. Matrix-assisted laser-desorption ionization-mass spectrometry imaging (MALDI-MSI) fulfills these constraints and several studies have indeed shown that this approach may characterize relevant tumor populations within cancer tissues [8–10]. In particular, because it relies on multivariate analysis of large-scale molecular signatures, MALDI-MSI has the potential to identify differences between histologically indistinguishable regions of a tissue, which are the basis of cancer phenotypes that drives tumor progression and ultimately determines which variant of cancer a patient will

experience. In addition, coupled with other MS analytical techniques, the approach cannot only define signatures but also identify characteristic biomolecular ions to get information on the underlining mechanisms of cancer cell resistance *in vivo* and eventually derive new targeted-treatments. Moreover, no requirement for specific staining/labelling agents is needed as compared with histological analysis, thus permitting multiplex analysis of several molecules in the same tissue section. Finally, it allows the correlation of molecular information with traditional histology by maintaining the spatial localization information of the analytes during mass spectrometric measurement.

Traditional analysis of MALDI-MSI images has made only limited use of the big amount of information achieved with this kind of data. To attain the distribution of molecular phenotypes while demonstrating the significant spatial heterogeneity present in the molecular phenotype distribution, the use of an efficient multivariate data processing pipeline is vital. The data analysis workflow proposed in this work is applied to fifteen tissues of experimental colorectal cancer images with different degree of heterogeneity using a mixture of isogenic cells all derived from HCT116 human colon adenocarcinoma cells. Using mixtures of Irinotecan-sensitive (S) and resistant (R) cell lines presenting highly similar molecular phenotypes will accurately represent the clinical question of clonal evolution and emergence of resistant cancer cells during treatment [11]. The strategy follows the three steps below:

a) Multiset multivariate curve resolution (MCR) on the 15 images. Basic MS spectral signatures and related distribution maps of all components in MSI are recovered. This method allows segmenting components related to background signal contributions from biological components. However, biological components related to different tissue types (necrotic and tumoral parts) may overlap, and hence no hard separation between them

could be achieved. Moreover, in this case, R and S cell lines cannot be easily separated by MCR due to highly similar MS signatures and/or unclear spatial pattern.

b) Multiset K-means image segmentation analysis on the MCR scores (concentration profiles) of only biological contributions. In this step, a hard clustering method separates classes of tumor and necrotic tissues. However, this unsupervised clustering is not sufficient/powerful enough to distinguish between R and S tumors, and hence, a supervised discrimination method is needed.

c) Partial-least squares discriminant analysis (PLS-DA) model on pixel spectra from tumoral clusters. Supervised classifier such PLS-DA using previous pixel labeling (classes R or S) will help to discriminate between irinotecan-resistant and irinotecan-sensitive cell lines even if only subtle differences among R and S are present.

The proposed approach is a method to find the spatial tumor heterogeneity based on MALDI-MSI measurements even in scenarios of highly similar cancer subpopulations.

2. Experimental

A model of tumor heterogeneity, using two human colon adenocarcinoma cell lines; one sensitive (HCT116) and another resistant (HCT116-SN50) to Irinotecan, was designed to examine the feasibility of the proposed strategy. The complete methodology and materials used are detailed in our recent paper [12]. Briefly, experimental tumors were generated in athymic mice by a subcutaneous injection of a mix of a resistant (R) and a sensitive (S) HCT116 cell lines at five ratios (100%S; 90%S; 50%S; 90%R; 100%R). Tumors were then collected, sliced, scanned with an Epson Perfection 4990 Photo Scanner and subsequently analyzed by a 4800 Plus MALDI TOF/TOFTM Analyzer. A total of 15 images corresponding to different tumor samples (corresponding to different mice) and replicates

of slices of the same tumor were analyzed. Table 1 summarizes the description and image labelling.

TABLE 1

3. Data analysis

The data analysis workflow proposed follow the four steps below:

- a) Image preprocessing
- b) Multiset MCR resolution on all images to recover basic spectral signatures and distribution maps of pure compounds contributions, allowing separation of tumor and necrotic contributions from background.
- c) Multiset K-means image segmentation analysis on the MCR scores (concentration profiles) of only biological contributions, to identify tumor or necrotic parts, permitting selection of tumor pixels for discrimination purposes.
- d) **Use of supervised classification method** to discriminate between irinotecanresistant and irinotecan-sensitive populations of tumor tissue.

These steps are described in detail in the following subsections.

3.1 Image preprocessing

To extract reliable conclusions and to retain maximum biological information from MSI data, it is necessary to choose appropriate preprocessing steps. We chose MALDIquant R package [13] to perform almost all preprocessing data, which includes baseline correction, spectra smoothing, peak detection and spectra alignment. A matrix **D** (*n*,*m*) of dimension *n* equal to ($x \times y$) pixels by *m* m/z values was generated per each image. Finally, a multiset

structure containing all 15 images was built, concatenating the 15 individual matrices one on top of each other. For a more detailed information about the preprocessing steps and parameters, please refer to our recent publication about the same dataset [12].

3.2 MCR-ALS resolution of MSI multiset

The goal of the image resolution step by Multivariate Curve Resolution Alternating Least Squares (MCR-ALS) algorithm is the decomposition of the data set into the distribution maps (relative amounts or concentration) and pure spectra of the signal contributions present in the imaged sample [14–16]. In this study, all 15 images were analyzed simultaneously in a single multiset structure. This multiset was obtained by appending the pixel MS matrix of the different MALDI images one on top of each other to form a column-wise augmented matrix, D_{aug} = [D100R1a; D100R1b; D100R1c; D100R2... D50R11]. This column-wise augmented data matrix D_{aug} can be decomposed using the bilinear model equation:

$$\mathbf{D}_{aug} = \mathbf{C}_{aug} \mathbf{S}^{T} + \mathbf{E}_{aug}$$
 Equation 1

where the matrix C_{aug} , is a column-wise matrix formed by as many submatrices C_i (matrix of the relative amounts or concentration of *n* components) as images in the multiset and S^T is a single matrix of pure spectra of *n* components, valid for all the images in the multiset. MCR-ALS estimates iteratively the matrices C_{aug} and S^T under the control of constraints

[12].

. In this study, the most common constraints in image resolution, such as non-negativity and normalization, were used. As a particularity, physiological information based on scanned images were coded as local rank constraints and were oriented to extract much more reliable MS signatures and distribution maps of the different tissue types present in the samples. Further information on the application of this local rank constraint can be found in [12]. Indeed, more details about the MCR-ALS method are given in [17–19] and a GUI to use the algorithm is freely available at <u>http://mcrals.info</u>.

MCR-ALS results, i.e., distribution maps (\mathbf{C} matrix) and pure spectra (\mathbf{S}^{T} matrix) may be further used to obtain additional information. In this work, distribution maps were used as starting information in K-means that allows a component selection of the profiles to be included in the segmentation process, i.e. resolved components linked clearly to background signal contributions are discarded and only those linked to biological contributions are retained for segmentation.

3.3 K-means multiset segmentation

The main goal of image segmentation is to partition the images into pixel groups built according to their pixel similarity in an unsupervised way [17]. In this work, K-means was used as a segmentation approach that tries to partition the dataset into K pre-defined clusters where each pixel belongs to only one group.

K-means analysis has been applied to the multiset structure containing all 15 images. MCR scores (concentration profiles) of biological contributions were used as starting information because MCR scores are compressed and noise-filtered and keep only biological relevant information on pixel composition, which makes them more suitable than raw spectra for segmentation purposes. It has been proved that dimensionality reduction prior to K-means clustering can be beneficial for the quality of the clustering [20].

Based on the comparison between segmentation map and grayscale images, we can determine which pixels are localized in each region and classify the corresponding spectra into categories (tumoral or necrotic tissues). Tumoral spectra will be used for further supervised classifications purposes.

3.4 Supervised Classification

Supervised learning methods are widely used in tissue MSI based-research when prior examples regarding píxel labeling in tissue samples is available [21]. In this work, Partial Least Squares Discriminant Analysis (PLS-DA) was used as supervised classification method to predict the class (R or S tumor type) of unknown samples according to their mass spectrum (MS pixels) [20]. The use of a supervised methodology responds to the fact that individual MS signatures of R and S cell subpopulations are not sufficiently different to be properly separated by the sole use of unsupervised methods [8]. PLS-DA is a discriminant technique based on PLS regression that builds a model using an appropriate set of latent variables to maximize the covariance between **X**, the MS pixels, and **Y**, a vector containing the categorical classes R and S, numerically coded. Since PLS-DA is very prone to provide overoptimistic results even in cross-validation [22], the use of permutation tests is also applied to investigate the statistical significance of the classifier accuracy [23].

In our case, labelling can only be done on images with 100% of either R or S cell types. Mixed images cannot be labelled since visually type R tumor cannot be differentiated from tumor type S. Therefore, from all 15 images only 8 images were labelled, please refer to table 1 to see which images correspond to 100% of one type of tumor. The calibration set comprises 3 images of 100% irinotecan-resistant (100R1a, 100Rab, 100R1c) from the same tumor and 3 images of 100% irinotecan-sensitive from the same tumor (100S3a, 100S3b, 100S3c). External validation set includes the images of 100% R and 100% S obtained from another grafted mouse but without replicates (100R2 and 100S4, respectively). Additionally, we have also applied the resultant predictive model to the

remaining images, which include mixtures of R and S tumors, in order to determine the spatial distribution of the tumors over the images.

Within the calibration test, cross-validation performed by leaving one image of each class out was used to choose the optimal number of latent variables as well as to test the predictive performance of the model, according to the cross-validation table shown below (see Table S1 in supplement material).

The classification rate is taken as a criterion of goodness for the developed model. The number of LV was selected as the one that provides a maximum CR in cross-validation. However, we applied the Parsimony principle: the inclusion of one more LV should represent a gain of more than 1 % in the CR value to be considered.

To asses the quality of the model, classification rate (accuracy), sensitivity (Se), defined as proportion of samples correctly classified within a particular class, selectivity (Sp), defined as proportion of samples correctly classified outside a particular class, and area under Receiver Operator Characteristic (AUROC), which is equal to 1 for perfect classification accross all possible values of the threshold. [24].

A permutation test has been used to assess whether the specific classification of the MS pixels in the two designed groups (R or S) is significantly better than random classification in two arbitrary groups. In this stduy, 6000 permutations have been done to define the classification parametres associated with the arbitrary class models. Statistical significance of the PLS-DA model is then assessed by relating the value of the classification parameters obtained by the original data set to the distribution of them calculated with the permuted data sets. In this way, permutation test allows to estimate a p-value for the classification rate obtained with the real labels. The significance threshold is usually set to

0.05 as in most biologic applications. p-values smaller than 0.05 indicate that differences between the two classes are statistically significant.

In order to understand the molecular ions that contribute to the discrimination between both classes (R and S tumor classes), the Variable Importance in Projection (VIP) has been used [25].

All data treatment (resolution, segmentation and classification methods) has been performed in Matlab (The MathWorks Inc.). The PLS-DA method has been applied using in-house routines, partly based on the PLS Toolbox (Eigenvector Research Inc.).

4. Results and discussion

4.1 Selection of tumoral pixels for further discrimination purposes.

Results from the proposed combination of MCR-ALS resolution with K-means segmentation in order to select the tumor spectra pixels of the tissues, that will be used for further discrimination tasks, could be found in our recent publication about the same dataset [12].

As an example, Figure 1a shows a schematic illustration of the proposed strategy for the image 100R1b. Firstly, basic description of pure contributions of the tissue is achieved by MCR-ALS. Background and biological (tumoral and necrotic) contributions could be clearly differentiated. Background signals present noise spectral signatures while biological components present chemically meaningful signatures. However, it was seen that distribution maps present overlaps among necrotic and tumoral contributions; i.e. no hard separation of these contributions is obtained and, due to heterogeneity, more than one contribution is needed to describe these two kinds of tissues [8]. Moreover, MCR results, i.e distribution maps (concentration profiles) and pure MS spectra (figure not shown), were not conclusive enough to separate R and S cell populations. Secondly, MCR concentration

profiles of biological contributions were used as input information for K-means segmentation. In this case, hard separation of necrotic and tumoral pixels is achieved. Two Clusters (green and orange) correspond to the tumor parts, they have the same spatial distribution of the dark grey regions in the scanned image. Similarly, two clusters (blue and brown) were associated with the necrotic parts. It is worth to mention that the information coming from segmentation maps is richer than that provided by the grayscale images because it reveals heterogeneity within necrotic and tumor tissues and is based on chemical information (from MS measurement) rather than simple color intensity. All the pixels could be assigned to any of the tissue types (tumoral or necrotic), as opposed to grayscale images, where the assignment of a pixel is just based on the grayscale level and hence some pixels cannot be straightforwardly assigned to any tissue type.

Now, we can determine which pixels are localized to each region and classify the corresponding spectra into categories (tumoral or necrotic tissues). However, this plain unsupervised clustering is not enough to distinguish between R and S tumors, and hence a dedicated method for discrimination is required. Therefore, tumoral spectra selected from the clusters associated with tumoral parts by K-means (see figure 1b) will be used as starting information in order to discriminate between R and S cell populations.

FIGURE 1

4.2. PLS-DA classification method. Discrimination between R and S populations

In supplementary material, Figure S1a shows the plots of the mean of CR and their confidence interval calculated from the binomial distribution in cross validation as a function of the increasing number of latent variables. The model seems to show a plateau

from 5 to 10 components and thus appears to be rather stable. Therefore, the number of LVs chose was 5 since inclusion of one more dimension did not represent a gain of more than 1% in the CR value.

Once the optimal number of LV has been selected, calculation the PLS-DA model on the calibration samples has been carried out. This calculation was performed selecting 5 LV and the same cross validation groups for internal validation. In Table S1 of supplementary material quality parameters obtained in the training set and in cross validation are collected. Similar classification performance was obtained in both cases; thus, we can consider the PLS-DA model reliable and stable.

High classification rate, CR, specificity and sensitivity were achieved. Both R and S pixel spectra were correctly classified into their corresponding class with more than 94% of accuracy in both calibration and cross-validation. Note that AUROC is nearly 1 in both cases and hence, practically perfect separation between the classes is achieved. The plot of sensitivity and specificity values as a function of the increasing class threshold for R class was presented in supplementary material (Figure S1b). Note that the same plot for S class would be complementary (figure not shown). The class threshold where the number of TP and TN is maximized, thus, better classification performance is achieved, corresponds to the point where the specificity line crosses the sensitivity line. It can be seen from this figure that the class threshold for the best classification of R class is ranged from 0.591 to 0.597 (0.403-0.409 for S class).

Although values of CR= 0.98 or AUROC = 0.99 could be considered vastly good and to correspond to a proper model with a high discriminating power, we have tested in those values can be reached purely by chance. In order to give a measure of the statistical significance of these quality parameters (p-value), a permutation test was carried out (see section 3.4). Statistical significance of CR and AUCROC of the PLS-DA model can be

evaluated by comparing them to the values of their null reference distributions H_0 obtained by permutation tests. Figure S2 of supplementary material shows the CR and AUCROC values obtained in calibration setup in red and for the permutations in gray. For both quality parameters there is a clear distinction between the null hypothesis distribution and the obtained estimation with the real labels. The results of the permutation test indicate that the specific classification is significant. If we compare the average value of the original classification with all the permutations, then a p-value can be obtained as described in section 3.4 (equations 5 and 6). Both p_{CR} (0.0002) and p_{AUCROC} (0.0002) values are smaller than 0.05, confirming the statistical significance of the obtained results.

The variables in the projection (VIPs) provided by the classification model were used to determine the most important ion molecular m/z values in the discrimination of the two classes. In Figure 2, the mean spectrum of pixels from R and S classes in black and red, respectively, have been overlapped with the VIPS to reveal the most relevant m/z values to distinguish both classes. It could be seen slight differences in spectral intensities which lead to understand some underlying behaviours of the different classes of images. Class R presents higher intensity at the m/z values of 291.1, 348.0, 672.0 and Class S present higher intensities at 520.3/522.3, 568.2, 738.5, 754.5/756.5, 770.5/772.5 and 798.5. Identification of these mass values could not be unequivocally done due to the resolution limits in the MS detection system. However, class R present higher intensities in masses around 750-800 that could assigned glycosyldiradylglycerols, be to glycerophosphocholines or glycerophosphoglycerols and to masses around 520-525 that could be related to oxidized glycerophospholipids or fatty acyl glycosides.

FIGURE 2

Once the model has been considered adequate, the true predictive performance of the classification model could be assessed by the external test set. Quality parameters from the test set confirm the same classification performance previously achieved by internal validation on the training samples (see external validation results in Table S2 in supplementary material). CR, specificity and AUCROC are similar to those obtained on the training set. However, lower value of sensitivity is obtained, consequently less ability to correctly recognize pixels belonging to R class is presented. It is worth to mention that no images from different mice of neither R nor S classes were included in the model development due to the limited number of images; Therefore, the variability coming from different mice was not considered in the model, explaining the slightly worse classification results for the external validation set. Nevertheless, since classification results obtained in external validation are satisfactory, the PLS-DA model could be considered adequate and its performance on future samples is expected to be comparable to those achieved on test samples.

Once the model is built and validated, it has been used to predict the distribution of R and S class pixels in the images coming from xenografts where both R and S cell lines in different proportions were inoculated. Figure 3 shows the pixel class assignment in the analyzed samples. As can be seen, pixels of the two cell lines are detected in tumors where both were inoculated. From the analyzed samples, no morphological distribution patterns can be associated with the growth of the two cell lines, i.e., both cell lines grow similarly and in a mixed way in the tissue analyzed. Qualitatively, it seems that both cell lines R and S present ratios in the tumor sections analyzed similar to the ratio inoculated, except for sample 90S8.

FIGURE 3

Indeed, this observation could be confirmed when quantitative ratios of R and S tumors were calculated by using the percentage of pixels predicted for each class in every sample over the total number of pixels considered to be tumor. An estimation of the percentage of R tumor populations is presented in Table 2.

TABLE 2

From Table 2, a satisfactory estimate of the percentage of R tumor populations can be observed except for the image 90S8. Without considering this image, a correlation of 0.93 between real % of R and estimated % of R was observed. This indicates that the % of the different cell lines found by PLS-DA in the analyzed images is closely correlated with the one in the inoculated samples. When grafted independently, the growth rate of the R and Scell lines are similar with a doubling time of 8.6 +/- 0.6 and 9.0 +/- 0.4 days respectively. The similarity in the predicted proportion of R cells between injection and resection suggests that there is no competition for growth between the two cell lines in vivo. The deviations in the % of cell lines predicted can be explained by biological variability between mice, but also because we analyzed a single tumor slice that cannot represent the composition of the full tumor. This is particularly evident in sample 90S8, where the results linked to the tissue section analyzed are far from the cell line R/S ratio expected.

As a final summary, the combination of MSI and a powerful protocol combining consecutive modeling steps of image resolution – segmentation – discrimination has proven to be efficient to distinguish the highly similar Rand irinotecan-sensitive tumor cell lines.

17

5. Conclusions

Combination of image resolution, segmentation and classification multivariate analysis method results in the capacity to much better discriminate highly similar tumors subpopulations (irinotecan-resistant (R) and irinotecan-sensitive (S)) in heterogenous cancer tissues investigated by MALDI-MSI.

In the present model of irinotecan-resistant human adenocarcinoma cells, discriminating differences in the intensities of the peaks of masses around 750-800 and 520-525 have been detected. These masses could be assigned to glycosyldiradylglycerols, glycerophosphocholines or glycerophosphoglycerols, and to oxidized glycerophospholipids or fatty acyl glycosides, respectively. Furthermore, the results obtained support the idea that the different tumor cell lines (R and S) present the same proliferation behavior *in vivo*, as shown by the preservation of ratios of R and S cells during tumor growth.

The conclusions reached in this work promote further biological research in order to evaluate the discriminant potential of this approach in case of bigger data sets showing a large variability among individuals. In particular, studies of patient samples are now required to more precisely characterize the type of heterogeneity that can be identified using this approach and the interest in patient monitoring during treatment [2,3]. In addition, a A systematic comparison with conventionally stained sections would be of interest to evaluate the potential of the proposed strategy as a complement to this routine approach.

Nevertheless, MALDI-MSI combined with the proposed data analysis protocol has been shown to be a valuable tool to investigate tumor heterogeneity even in scenarios of highly similar cancer subpopulations.

18

6. Acknowledgements

This work is part of the BEST Postdoctoral Program, funded by the European Commission under Horizon 2020 Marie Skłodowska-Curie Actions COFUND scheme (Grant Agreement no. 712754) and by the Severo Ochoa program of the Spanish Ministry of Science and Competitiveness (Grant SEV-2014-0425 (2015-2019)). A.J. acknowledges financial support from the Catalan government through project 2017 SGR 753 and the Spanish government through project CTQ2015-66254-C2-2-P. We would like to acknowledge the Departament d'Universitats, Recerca i Societat de la Informació de la Generalitat de Catalunya (expedient 2017 SGR 1721); the Comissionat per a Universitats i Recerca del DIUE de la Generalitat de Catalunya; and the European Social Fund (ESF). Additional financial support has been provided by the Institut de Bioenginyeria de Catalunya (IBEC). IBEC is a member of the CERCA Programme/Generalitat de Catalunya. This publication has been also funded with support from the French National Research Agency under the program "Investissements d'avenir" Grant Agreement LabEx MAbImprove: ANR-10-LABX-53. A. T. and P. M. acknowledge the support of the Fondation pour la Recherche Médicale. We thank the Experimental Histology Network of Montpellier for histology/immunohistology experiments (RHEM, http://www.rhem.cnrs.fr).

7. Conflict of interest

The authors report there are not conflicts of interest

6. Bibliography

 F.A. Haggar, R.P. Boushey, D. Ph, Colorectal Cancer Epidemiology : Incidence, Mortality, Survival, and Risk Factors, 6 (2009) 191–197. doi:10.1055/s-0029-1242458.

- I. Dagogo-Jack, A.T. Shaw, Tumour heterogeneity and resistance to cancer therapies, Nat. Rev. Clin. Oncol. 15 (2017) 81. https://doi.org/10.1038/nrclinonc.2017.166.
- [3] N. McGranahan, C. Swanton, Clonal Heterogeneity and Tumor Evolution: Past,Present, and the Future, Cell. 168 (2017) 613–628. doi:10.1016/j.cell.2017.01.018.
- M. Jamal-Hanjani, G.A. Wilson, N. McGranahan, N.J. Birkbak, T.B.K. Watkins, S. [4] Veeriah, S. Shafi, D.H. Johnson, R. Mitter, R. Rosenthal, M. Salm, S. Horswell, M. Escudero, N. Matthews, A. Rowan, T. Chambers, D.A. Moore, S. Turajlic, H. Xu, S.-M. Lee, M.D. Forster, T. Ahmad, C.T. Hiley, C. Abbosh, M. Falzon, E. Borg, T. Marafioti, D. Lawrence, M. Hayward, S. Kolvekar, N. Panagiotopoulos, S.M. Janes, R. Thakrar, A. Ahmed, F. Blackhall, Y. Summers, R. Shah, L. Joseph, A.M. Quinn, P.A. Crosbie, B. Naidu, G. Middleton, G. Langman, S. Trotter, M. Nicolson, H. Remmen, K. Kerr, M. Chetty, L. Gomersall, D.A. Fennell, A. Nakas, S. Rathinam, G. Anand, S. Khan, P. Russell, V. Ezhil, B. Ismail, M. Irvin-Sellers, V. Prakash, J.F. Lester, M. Kornaszewska, R. Attanoos, H. Adams, H. Davies, S. Dentro, P. Taniere, B. O'Sullivan, H.L. Lowe, J.A. Hartley, N. Iles, H. Bell, Y. Ngai, J.A. Shaw, J. Herrero, Z. Szallasi, R.F. Schwarz, A. Stewart, S.A. Quezada, J. Le Quesne, P. Van Loo, C. Dive, A. Hackshaw, C. Swanton, Tracking the Evolution of Non–Small-Cell Lung Cancer, N. Engl. J. Med. 376 (2017) 2109–2121. doi:10.1056/NEJMoa1616288.
- [5] P.C. Nowell, The clonal evolution of tumor cell populations, Science (80-.). 194
 (1976) 23 LP 28. doi:10.1126/science.959840.
- [6] Z. Piotrowska, M.J. Niederst, C.A. Karlovich, H.A. Wakelee, J.W. Neal, M. Mino-Kenudson, L. Fulton, A.N. Hata, E.L. Lockerman, A. Kalsy, S. Digumarthy, A.

Muzikansky, M. Raponi, A.R. Garcia, H.E. Mulvey, M.K. Parks, R.H. DiCecca, D. Dias-Santagata, A.J. Iafrate, A.T. Shaw, A.R. Allen, J.A. Engelman, L. V Sequist, Heterogeneity Underlies the Emergence of & the emergence of & the emergence of the emergence of the terms and the terms and the emergence of the terms and terms

- [7] L.A. Diaz Jr, R.T. Williams, J. Wu, I. Kinde, J.R. Hecht, J. Berlin, B. Allen, I. Bozic, J.G. Reiter, M.A. Nowak, K.W. Kinzler, K.S. Oliner, B. Vogelstein, The molecular evolution of acquired resistance to targeted EGFR blockade in colorectal cancers, Nature. 486 (2012) 537–540. doi:10.1038/nature11219.
- [8] S.M. Willems, A. Van Remoortere, Imaging mass spectrometry of myxoid sarcomas identifies proteins and lipids specific to tumour type and grade, and reveals biochemical intratumour heterogeneity §, (2010) 400–409.
- R. Emrys A. Jones, Alexandra van Remoortere, L.A.M. ´J. M. van Zeijl, Pancras C.
 W. Hogendoorn, Judith V. M. G. Bovée, André M. Deelder, Multiple Statistical Analysis Techniques Corroborate Intratumor Heterogeneity in Imaging Mass Spectrometry Datasets of Myxofibrosarcoma, 6 (2011).
 doi:10.1371/journal.pone.0024913.
- B. Balluff, C.K. Frese, S.K. Maier, C. Schöne, B. Kuster, M. Schmitt, M. Aubele, H. Höfler, A.M. Deelder, A.J.R. Heck, P.C.W. Hogendoorn, A.F.M. Altelaar, A. Walch, L.A. Mcdonnell, De novo discovery of phenotypic intratumour heterogeneity using imaging mass spectrometry, (2015) 3–13. doi:10.1002/path.4436.

- [11] M. Greaves, C.C. Maley, Clonal evolution in cancer, Nature. 481 (2012) 306–313. doi:10.1038/nature10762.
- [12] S. Mas, A. Torro, N. Bec, L. Fernández, G. Erschov, C. Gongora, C. Larroque, P. Martineau, A. de Juan, S. Marco, Use of physiological information based on grayscale images to improve mass spectrometry imaging data analysis from biological tissues, Anal. Chim. Acta. (2019). doi:10.1016/j.aca.2019.04.074.
- S. Gibb, K. Strimmer, Maldiquant: A versatile R package for the analysis of mass spectrometry data, Bioinformatics. 28 (2012) 2270–2271.
 doi:10.1093/bioinformatics/bts447.
- [14] J. Jaumot, R. Tauler, Potential use of multivariate curve resolution for the analysis of mass spectrometry images, Analyst. 140 (2015). doi:10.1039/c4an00801d.
- [15] C. Bedia, R. Tauler, J. Jaumot, Analysis of multiple mass spectrometry images from different Phaseolus vulgaris samples by multivariate curve resolution, Talanta. 175
 (2017) 557–565. doi:10.1016/j.talanta.2017.07.087.
- S. Piqueras, C. Krafft, C. Beleites, K. Egodage, F. von Eggeling, O. Guntinas-Lichius, J. Popp, R. Tauler, A. de Juan, Combining multiset resolution and segmentation for hyperspectral image analysis of biological tissues, Anal. Chim. Acta. 881 (2015) 24–36. doi:10.1016/j.aca.2015.04.053.
- [17] R. de Juan, Anna; Rutan, S. and Tauler, Two-Way Data Analysis: Multivariate
 Curve Resolution Iterative Resolution Methods, in: B. Brown, S. D.; Tauler, R.
 and Walczak (Ed.), Compr. Chemom., Elsevier, 2010: pp. 325–344.
- [18] J. Jaumot, A. de Juan, R. Tauler, MCR-ALS GUI 2.0: New features and applications, Chemom. Intell. Lab. Syst. 140 (2015) 1–12.

22

doi:10.1016/j.chemolab.2014.10.003.

- [19] R. Tauler, Multivariate curve resolution applied to second order data, Chemom. Intell. Lab. Syst. 30 (1995) 133–146. doi:10.1016/0169-7439(95)00047-X.
- [20] W. Liu, K. Yuan, D. Ye, Reducing microarray data via nonnegative matrix factorization for visualization and clustering analysis, J. Biomed. Inform. 41 (2008) 602–606. doi:10.1016/j.jbi.2007.12.003.
- [21] Y. Zhang, X. Liu, Machine learning techniques for mass spectrometry imaging data analysis and applications, Bioanalysis. (2018) 10–13. doi:10.4155/bio-2017-0281.
- [22] R. Rodríguez-Pérez, L. Fernández, S. Marco, Overoptimism in cross-validation when using partial least squares-discriminant analysis for omics data: a systematic study, Anal. Bioanal. Chem. 410 (2018) 5981–5992. doi:10.1007/s00216-018-1217-1.
- [23] M. Ojala, Permutation Tests for Studying Classi er Performance, J. Mach. Learn. Res. 11 (2009) 1833–1863. doi:10.1109/ICDM.2009.108.
- [24] H. Abdi, Signal Detection Theory (SDT), (1966) 1–9.
- [25] S. Wold, E. Johansson, M. Cocchi, PLS: Partial Least Squares Projections to Latent Structures, in: 3D QSAR Drug Des. Vol. 1 Theory Methods Appl., 1993.

Table 1. Percentage of both R and S cell lines in xenografts, replicate number of either thesame or different tumors (mice) and code of the image sections analyzed.

Percentage of cell lines	Replicate of the same tumor	Replicate number of different tumors	Image code
100% R	a	1	100R1a
100% R	b	1	100R1b
100% R	с	1	100R1c
100% R	-	2	100R2
100% S	a	1	100S3a
100% S	b	1	100S3b
100% S	С	1	100S3c
100% S	-	2	100S4
90% R	-	1	90R5
90% R	-	2	90R6
90% S	-	1	90S7
90% S	-	2	90S8
90% S	-	3	90S9
50% S	-	1	50R10
50% S	-	2	50R11

Table 2. Real vs. estimated % of R cell line in heterogenous tumor samples as predicted inthe analyzed MS images by PLS-DA.

Image	Real % of R	Estimated % of R
90R5	90	79
90R6	90	91
90\$7	10	14
90\$8	10	74
90\$9	10	8
50R10	50	39
50R11	50	76

Figures



Figure 1. (a) Graphical representation of the strategy used for PLS-DA-oriented pixel selection for the 100R1b image. (b) grayscale images (left side) and selected tumoral pixels for further classification tasks (right side) for all images.



Figure 2. Mean spectrum of pixels from images of R class (black), mean spectrum of the MS pixels from images of S class (red) and VIP values higher than 10 (green).VIP are normalized and spectra dived in 3 ranges of masses (250-450, 450-650, 650-900) for a better visualization. Most relevant m/z values for discrimination among R and S class have been displayed.



Figure 3. Distribution of R and S tumor populations in the heterogenous images.