

# Statistical treatment of 2D NMR COSY spectra in metabolomics: data preparation, clustering-based evaluation of the Metabolomic Informative Content and comparison with $^1\text{H}$ -NMR

Baptiste Féraud<sup>1,2</sup> · Bernadette Govaerts<sup>1</sup> · Michel Verleysen<sup>2,3</sup> · Pascal de Tullio<sup>4</sup>

Received: 3 October 2014 / Accepted: 13 July 2015  
© Springer Science+Business Media New York 2015

**Abstract** Compared with the widely used  $^1\text{H}$ -NMR spectroscopy, two-dimensional NMR experiments provide more sophisticated spectra which should facilitate the identification of relevant spectral zones or biomarkers in metabolomics. This paper focuses on  $^1\text{H}$ - $^1\text{H}$  COReLation SpectroscopY (COSY) spectral data. In spite of longer inherent acquisition times, it is commonly accepted by users (biologists, healthcare professionals) that the introduction of an additional dimension probably represents a huge qualitative step for investigations in terms of metabolites identification. Moreover, it seems natural that more information leads to more predictive power. But, until now, very few statistical studies clearly proved this assumption. Therefore a fundamental question is “Is this supplementary information relevant?”. In order to extend the statistical properties developed for 1D spectroscopy to the challenges raised by 2D spectra, a rigorous study of the performances of COSY spectra is needed as a prerequisite. Having introduced new pre-processing concepts, such as the Global Peak List or an ad hoc 2D “bucketing”, this paper presents an innovative methodology based on

multivariate clustering algorithms to evaluate this question. Numerical clustering quality indexes and graphical results are proposed, based both on the spectral presence or absence of peaks (binary position vectors) and on peak intensities, and through different levels of spectral resolution. The second goal of this paper is to compare clustering performances obtained on COSY and on  $^1\text{H}$ -NMR spectra, with the aim of understanding to what extent the COSY spectra carry more Metabolomic Informative Content about the signal than 1D ones. The methodology is applied to two real experimental designs involving different groups of spectra (which define the signal): a 4-mixture cell culture media containing various supervised metabolites and a complex human serum based design. It is shown that COSY spectra appear to be statistically powerful and, in addition, provide better clustering results than corresponding  $^1\text{H}$ -NMR when using unlabeled information. Consequently, additional information appears to be relevant for metabolomics applications.

**Keywords** COSY spectra · Metabolomic Informative Content · Peak lists · Multivariate clustering algorithms · Real data

✉ Baptiste Féraud  
baptiste.feraud@uclouvain.be

<sup>1</sup> Institute of Statistics, Biostatistics and Actuarial Sciences (ISBA), Université Catholique de Louvain (UCL), Voie du Roman Pays 20, bte L1.04.01, B-1348 Louvain-la-Neuve, Belgium

<sup>2</sup> Machine Learning Group, Université Catholique de Louvain (UCL), Voie du Roman Pays 20, bte L1.04.01, B-1348 Louvain-la-Neuve, Belgium

<sup>3</sup> SAMM, Université Paris I, Panthéon - Sorbonne, Paris, France

<sup>4</sup> Center for Interdisciplinary Research on Medicines (CIRM), Université de Liège (ULg), Liège, Belgium

## 1 Introduction

Widely used for finding bio-active metabolites in living organisms and to study the impact of various biotic and abiotic effects, metabolomics studies (Nicholson et al. 2002) were initially intended to search for metabolites interpretable as “biomarkers” for a given disease or toxicity. These discriminant metabolites are usually tracked in partial or complete spectral images, linked with biofluids (serum, urine) or tissues, using adequate univariate or

multivariate statistical techniques. It is clear that the sampling preparation, spectral acquisition and data pre-processing have a crucial impact on the global quality of the data and on the chances to finally discover relevant biomarkers. Furthermore, an accurate statistical analysis will rarely compensate “dirty” initial data. Many acquisition and pre-processing tools are available and it is often difficult to objectively and quantitatively evaluate which ones will provide the more informative data in a given context. This motivated the development of the methodology described in this paper when data are organized in groups.

In this context, proton nuclear magnetic resonance ( $^1\text{H}$ -NMR) spectroscopy generates spectral profiles describing the metabolite composition of collected samples. A comparison of several spectra of metabolites in various specific states permits a preliminary graphical and qualitative investigation of changes in metabolite composition inherent to the presence of a “stressor”. However, the complexity of  $^1\text{H}$ -NMR spectra and the number of spectra (of samples) usually available in metabolomics studies require a semi-automated data analysis. In addition, systematic differences between samples are often hidden behind biological noise and/or behind peak shifts.

To avoid these peak shifts and the usual overlappings of potentially independent signals, the use of two-dimensional homonuclear or heteronuclear experiments has gained attention these last years to identify metabolites. However, these tools are often left behind due to longer, and sometimes prohibitive, acquisition times. In this paper, COSY spectra are investigated (Aue et al. 1976; Keeler 2010; Akitt and Mann 2000). COSY consists in a correlation-based method for determining which signals arise from neighboring protons (usually up to four bonds). Correlations appear when there is spin–spin coupling between protons (i.e. correlation between two or more nearby chemical processes). Several works exist in the literature in order to get around the time-consuming acquisition problem, specifically for multidimensional COSY data [for example Ultrafast-COSY (Giraudeau et al. 2009; Queiroz et al. 2013) or PALS-COSY (Vega et al. 2010)]. Compared to  $^1\text{H}$ -NMR experiments for instance, the introduction of an additional dimension should allow a better representation of metabolites, a better predictive power and a better biomarker identification. Several works are already based on 2D-COSY or on faster extensions of COSY, such as in Xi et al. (2007), Yun et al. (2013) or Le Guennec et al. (2014). But, until now, very few statistical studies clearly tackled the potential superiority of 2D spectra [these studies mainly concern the group of Bruschweiler, see for instance Bruschweiler and Bingol (2011) and Bruschweiler et al. (2014)]. This lack leads to a central and

fundamental question: is this supplementary information really relevant in the context of metabolomics analyses? If it is the case, the interest of using two-dimensional methods would be legitimated, even at the price of potentially superior acquisition times.

In this paper, the innovative notion of “Metabolomic Informative Content” (MIC) is central. The intention is to systematically evaluate the amount of captured information (i.e. the extent to which signals are captured) compared with the noisy part through several spectra. It goes away from the more classic version of repeatability in spectrometry which involves only one spectrum at a time. Generally, a measure can be considered as the sum of a signal (useful information) and noise. In many metabolomics studies, the signal is controlled and often linked with the existence of different mixtures, different samples or different groups of people (people affected by a disease vs. healthy people, people affected by a disease at different levels, ...). And the noisy part of the information is generally provoked by the characteristics of the design, the methods of acquisition, the temporal dimension (repetitions during time, deterioration of the samples, etc...) and the processing. The purpose would be to measure the predominance of the signal with regard to noise, but it is not obvious in a non-supervised context. A way is to suppose that only the useful information can be “repeatable” and that the noise is independent and identically distributed, with  $\text{signal} \perp \text{noise}$ . With this hypothesis, evaluating the Metabolomic Informative Content (MIC) gives an idea if  $\text{signal} \gg \text{noise}$  or not. In other words, for any set of spectra, it allows a direct comparison between different spectral tools and can also be seen as an evaluation of predictive ability in this paper.

According to this context and by using peak lists data sets, the first goal is to show that COSY spectra are successful, i.e. that COSY spectra allow to capture the main part of the information connected to the signal(s). The second goal is to demonstrate that COSY spectra are “more successful” than  $^1\text{H}$ -NMR corresponding spectra and, by doing so, to demonstrate the utility and the importance of the additional second dimension. To reach these goals and to obtain quantitative responses, a multivariate clustering approach is selected and applied on unlabeled spectra in order to recover the signal. Two clustering algorithms (hierarchical Ward algorithm and K-Means algorithm) are used, as well as multiple combinations between different distance measures, between the use of binary positions vectors (presence or absence of a peak) and of intensity vectors and between different spectral resolutions.

The paper is organized as follows. Section 2 provides a detailed description of the two experimental designs used to illustrate the methodology and reach the goals. These

experimental designs are both involving real data: the first one implies repeated measurements of 4-mixture cell culture media containing various supervised metabolites, and the second one, a human serum based design with time sampling repetitions and multiple measurement permutations. In each case, COSY spectra and  $^1\text{H}$ -NMR spectra are collected together for further comparisons. Section 3 provides a description of the different pre-processing steps: construction of a “Global Peak List”, construction of binary positions vectors and intensity vectors, symmetrisation, water peak deletion, outlier deletion, etc... In this section, it is also explained how the spectral resolution is controlled (and the size of the data sets) by performing a so-called “bucketing” based on changes in the number of decimals in the global matrix. Clustering algorithms and associated distance measures are also detailed in this section. Section 4 contains the results: first, an analysis of the COSY spectra performances based on both positions (presences or absences of peaks) and intensities. Numerical outputs and a dendrogram illustrate these results. The comparisons between 1D and 2D spectra are also discussed in Sect. 4. Finally, a general conclusion and further works are given in Sect. 5.

## 2 Materials: COSY spectra, experimental designs and acquisition parameters

In this section, the second dimension in COSY spectra will motivate the main goal of this paper. Details about the two real experimental designs used to illustrate the methodology are provided. And finally, a technical explanation of  $^1\text{H}$ -NMR and COSY spectra acquisition is proposed in Sect. 2.2.

### 2.1 Experimental designs

The methodology presented in this paper is applied on two real data sets, obtained from two designed experiments.

#### 2.1.1 First design: cell culture media

The first design is based on four different mixtures (four cell culture media containing various levels of different metabolites like fetal bovine serum, amino acids, vitamins, proteins, etc...): DMEM/F12, MEM, RPMI/1640 and DMEM. Five hundred micro litre of cell culture media were supplemented with 200  $\mu\text{l}$  of deuterated phosphate buffer (0.1 M, pH=7.4) containing 1 % of sodium azide and 10  $\mu\text{l}$  of TMSP (10 mg/ml). The solutions were then transferred into 5 mm NMR tubes before NMR measurement. Three samples per mixture were collected, and three repeated measures were performed on each sample. All

samples were subject to freezing ( $-20^\circ\text{C}$ ) and defrosting steps before the analysis, with real risks of degradation and bacterial contamination because of the duration of the 2D analysis process.

In this design, signal corresponds to the four initial mixtures and noise arises from the sampling, time replicates, risks of degradation and other acquisition and condition parameters. Thirty six measures are finally available, corresponding to 36 COSY spectra and 36 corresponding peak lists. These peak lists are  $(t_i \times 3)$  matrices: with  $t_i$  the number of detectable peaks (whose values are above a machine-designed threshold) in sample  $i$  ( $i = 1, \dots, 36$ ). The three columns in these matrices correspond to the coordinate, or chemical shift, on the first axis (ppm), the coordinate on the second axis (ppm) and the raw intensity of the peak. The 36 corresponding  $^1\text{H}$ -NMR spectra were also collected. The total number of spectra is called  $n$  in the remainder of the paper.

#### 2.1.2 Second design: human serum

The second experimental design is based on human serum. Peripheral blood was collected in serum separating tubes (Greiner). Sera were distributed into 0.5 ml aliquots and stored at  $-80^\circ\text{C}$  just after sampling. Four blood donors were engaged for the study. For each collected sample, 500  $\mu\text{l}$  of serum, defrosted just before the analysis, were supplemented with 200  $\mu\text{l}$  of deuterated phosphate buffer (0.1 M, pH = 7.4) containing 1 % of sodium azide and 30  $\mu\text{l}$  of TMSP (10 mg/ml). The solutions were then transferred into 5 mm NMR tubes before NMR measurement. The design consists in eight days of measurements with replicates within each day and multiple permutations according to a latin hypercube sampling (LHS) method (Iman 2008) (in order to avoid confusion between donors and times of analysis). For each day, the four donors samples were analyzed and provided four COSY spectra and four  $^1\text{H}$ -NMR spectra. Spectral techniques (1D or 2D COSY) have not been applied at the same moment of the day, thus creating different delays before spectral measurement.

Finally, eight measures/spectra/peak lists are obtained per donor, which corresponds to 32 measures in all (32 1D spectra and 32 COSY spectra). Again, peak lists are from particular interest and correspond to  $(t_i \times 3)$  matrices with  $t_i$  the number of detectable peaks in the sample, or spectra,  $i$  ( $i = 1, \dots, 32$ ).

### 2.2 Spectra acquisition

1D and 2D spectra were recorded at 298 K on a Bruker Avance spectrometer operating at 500.13 MHz for the proton signal acquisition. The instrument was equipped

with 3 channels and with a 5mm triple resonance (HCN) cryoprobe with a Z-gradient and automatic tuning and matching system (ATMA). All samples were locked using deuterated water and shims were tuned using auto-tune (Topshims) and samples are measured manually without spinning. Due to the nature of the samples, a presaturation sequence was used in all the experiments in order to minimize the water signal. All data were referenced to internal sodium 3-trimethylsilyl-2,2,3,3-d<sub>4</sub>-propionate (TMSP) at 0.00 ppm chemical shift (all spectra are calibrated with regard to TMSP).

According to sample type, the <sup>1</sup>H-NMR spectra were acquired using either a 1D NOESY-presat sequence (cell culture mixture) or a CPMG relaxation-editing sequence with presaturation (human sera). Upon the presence of proteins in serum, the use of a sequence with a T2 filter (CPMG) greatly improves the baseline.

The NOESY-presat experiment (pulse sequence `noesygppr1d` supplied by Bruker) used a RD-90°-t1-90°-tm-90°-sequence with a relaxation delay of 4 s, a mixing time of 100 ms and a fixed t1 delay of 20 μs (spectral window of 10000 Hz). The water suppression pulse was placed during the relaxation delay (RD). The number of transients was typically 32. The acquisition time was fixed to 3.2769001s and a quantity of four dummy scans was chosen. FIDs were collected in 64 K time data points. The data were processed with the Bruker Topspin 2.1 software with a standard parameter set. The phase and baseline corrections were performed manually over the entire spectral range and a line broadening of 0.3 Hz was applied.

The CPMG experiment (pulse sequence `cpmgrp1d` supplied by Bruker) used a RD-90°-(t-180°-t<sub>n</sub>)-sequence with a relaxation delay (RD) of 2s, a spin echo delay (t) of 400 μs and the number of loops (n) equal to 80. The water suppression pulse was placed during the relaxation delay (RD) and a spectral window of 10245 Hz. The number of transients was typically 32. The acquisition time was fixed to 3.982555 s and a quantity of four dummy scans was chosen. FIDs were collected in 64 K time data points. The data were processed with the Bruker Topspin 2.1 software with a standard parameter set. The phase and baseline corrections were performed manually over the entire spectral range without line broadening.

Gradient enhanced magnitude COSY experiment (pulse sequence `cosygppr1d` supplied by Bruker) with a presaturation during relaxation delay was used for 2D measurements. Spectra were collected with 4096 points in t2 and 300 points in t1 over a sweep width of 10 ppm, with six scans per t1 value. The acquisition times were fixed to 0.2557028s in F2 and 0.0187212 in F1. The resulting COSY spectra were processed in Topspin 2.1 using standard methods, with sine-squared apodization in both

dimensions and zero filling in t1 to yield a transformed 2D dataset of 2048 by 2048 points.

Individual peak lists were then extracted using ACD/Labs 12.00 (ACD/NMR processor, freeware). For each 2D spectra, a ( $t_i \times 3$ ) matrix provides the ( $X, Y$ ) coordinates of the  $t_i$  peaks observed in the spectrum and their related intensities in a third column.

### 3 Methods: global peak list, pre-processing steps and clustering analysis

In this section, all data manipulations and pre-processing steps, both for standard metabolomics studies and for the particular purpose of this paper, are detailed. Then, a discussion on how to control the data resolution and, consequently, the size of the databases is proposed. Finally, clustering algorithms and related distance measures are also described in detail.

#### 3.1 Global peak list matrix

When one wants to perform simultaneous data analysis of a set of COSY spectra expressed in peak lists, a first requirement is to gather them together in a global object. The solution proposed here consists in building a so-called “Global Peak List” (GPL) matrix from the  $n$  ( $t_i \times 3$ ) individual peak list matrices available for each of the  $n$  2D spectra. This ( $T \times N$ ) GPL matrix includes, as rows, the  $T$  pairs of coordinates that appear in at least one of the individual spectra. The  $N = 2 + 2n$  columns include first the two peak coordinates columns, and then, for each individual spectrum, two columns corresponding to the observed peak intensities and to a deduced binary number (1 or 0) indicating if the peak appears or not in the spectrum (i.e. corresponds then to a strictly positive or to a null intensity). For the first design, the GPL is a ( $3250 \times 74$ ) matrix ( $74 = 2 + 2 \times 36$ ). For the second one, the GPL is a ( $6686 \times 66$ ) matrix ( $66 = 2 + 2 \times 32$ ). The GPL matrix can also be viewed as a combination of three matrices: a ( $T \times 2$ ) matrix of coordinates, a ( $T \times n$ ) matrix of intensities  $I$  and a ( $T \times n$ ) matrix of positions  $P$  (presence or absence).

In this paper, both spectral intensities and peak positions are considered of particular interest. Working on the signals’ positions or, in other words, on the simple existence of signals is motivated by a biological justification: a signal, or a particular metabolite, can be observed or not for a particular donor or in a particular media. If a signal is present in detectable quantity, this presence or absence is supposed to be stable, whereas intensities are variable from one measure to another, according to potential uncontrolled

factors. Working on positions, or absence/presence, can then be seen as a qualitative approach; working on intensities can be seen as a more quantitative approach as it is directly linked with concentrations.

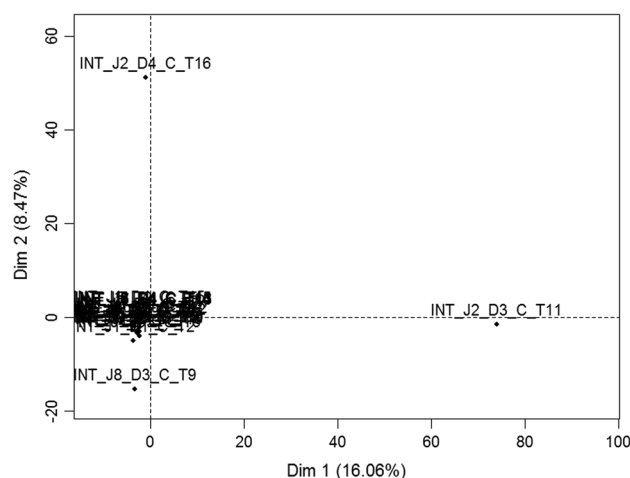
### 3.2 Pre-processing steps

The following pre-processing steps have been implemented on the GPL matrices:

- *Symmetrisation* By construction, all COSY spectra have to be symmetrical around the diagonal (see Sect. 2.1). Consequently, all other points or peaks are artifacts and have to be removed. It mainly concerns negative intensities and typical “crosses” which arise when choosing a wrong baseline and/or when some signals are abnormally too intense (Marion and Bax 1988).
- *Water zone deletion* Water is not of interest for metabolomics investigations. In COSY spectra, water peaks are concentrated in a square area between 4.5 and 5.5 ppm. Intensities and positions inside this area were simply removed by assigning them a zero value. A light loss of information can occur concerning metabolites present in this zone. More advanced methods also exist, see for example Mao and Ye (1997).
- *Normalization of the intensities* Vectors of recoded intensities are simply obtained by applying a constant sum ( $CS = 1$ ) normalization (Craig et al. 2006) after water zone deletion. This method is also known as 1-Norm (Rasmussen et al. 2011).
- *Outlier deletion* Some spectra can contain unexpected extreme signal values due, for example, to exceptional external factors. These outliers should be identified before data analysis and removed because they could be too influent. On both position and recoded intensities vectors, principal component analysis (PCA) were applied for graphical identification of outliers, and Mahalanobis distances were calculated between raw spectra and between group centered spectra. In the context of the first experimental design, three spectral outliers were found and ignored for upcoming analysis (because simultaneously suspected via positions and intensities). In the same way, three spectra were suspected to be outliers for the second design (see Fig. 1 for an illustration).

### 3.3 Control of data resolution

In  $^1\text{H}$ -NMR spectroscopy, bucketing tools are common and widely used to control the spectral resolution and/or to overcome the misalignments problem. Classical and more advanced bucketing methods have already shown their



**Fig. 1** Graphical identification of outliers using PCA's score plot. Here, three spectral outliers are detected based on intensities vectors: Day 2 (J2) at Time 16 (T16) for donor 4 (D4), Day 2 (J2) at Time 11 (T11) and Day 8 (J8) at Time 9 (T9) for donor 3 (D3)

usefulness for  $^1\text{H}$ -NMR spectra (Rousseau 2011; Sousa et al. 2013). In this paper, a bucketing step adapted to 2D-COSY is proposed in order to control the size of the database and, consequently, the resolution level of the two-dimensional spectra. Practically, a variation of the number of decimals of the coordinates is proposed. The intensities belonging to a bucket are then aggregated. For example, if the couples of coordinates [3.286; 4.194], [3.281; 4.189] and [3.278; 4.191] provide positive intensities INT1, INT2 and INT3 respectively, the couple [3.28; 4.19] provides an intensity equal to  $\text{INT1} + \text{INT2} + \text{INT3}$  when adjusting the number of allowed decimals from 3 to 2. In this paper, three, two and one decimal cases are tested for analysis. Using this method, the width of the COSY peaks is adjusted and the resolution and size of the databases are also adjusted simultaneously. Furthermore, intermediate resolutions can of course be computed in the same way.

For positions, the aggregation process leads to another dummy variable and results from a simple function: a peak is considered in a bucket if at least one peak is present at the lower resolution level. For example, 1 and 1 lead to 1, 1 and 0 lead to 1 and 0 and 0 lead to 0 (absence everywhere). Finally, for the two experimental designs, the GPL matrices have the dimensions described in Table 1 after “bucketing” and before outlier deletion.

**Table 1** Dimensions of the GPL matrices according to different resolutions

	One decimal	Two decimals	Three decimals
First design	(909 × 74)	(2348 × 74)	(3250 × 74)
Second design	(1106 × 66)	(4172 × 66)	(6686 × 66)



### 3.4 Clustering algorithms and distance measures

This section introduces a methodology able to quantify and compare with adequate indexes the MIC (in a sense of advantageous signal capture compared to noise) of different sets of spectra. This methodology is applied in Section 4 to the two designs (each organized in four groups: mixtures or blood donors) in order to compare the MIC of 1D and 2D COSY spectra.

For 2D experiments, several combinations between intensities and position vectors, and between the use of one, two or three decimal(s) GPL matrices is tested. Using all this information, an intuitive way to evaluate the MIC of COSY spectra consists in non-supervised multivariate clustering (blind, with no a priori labeled information). The key idea is quite simple: if one manages to well separate and recover the four initial mixtures starting from the 36 unlabeled spectral measures of the first design, and/or if one manages to well separate and recover the four blood donors starting from the 32 unlabeled spectral measures of the second design, the goal is reached. It would mean that the signal, specific to each group of spectra, is sufficiently captured, despite the noise due to the time repetitions, the sampling methods, the freezing and defrosting periods of the samples, the risks of bacterial contamination, the potential environmental changes, etc...In other words, the objective is to verify that the within-cluster (intra) variance is minimized and that the between-cluster (inter) variance is maximized during data acquisition. Clustering is a natural and accurate tool to check this problem.

Two well-known clustering algorithms are considered: the hierarchical Ward algorithm and the K-Means algorithm (for intensities only in this second case).

- *The Ward's algorithm* (Ward 1963; Murtagh and Legendre 2011) is a commonly used procedure for forming hierarchical groups of mutually exclusive subsets. It is particularly useful for large-scale studies (ideally with more than 100 objects) when a precise optimal solution for a specified number of groups is not practical. Given  $n$  objects, this procedure reduces them to  $n - 1$  mutually exclusive sets by considering the union of all possible  $n(n - 1)/2$  pairs and selecting the union having the minimal dissimilarity measure or aggregation index. Initially, each object is assigned to its own cluster and then the algorithm proceeds iteratively, at each stage joining the two most similar clusters, continuing until there is just a single cluster. At each stage distances between clusters are recomputed by the Lance–Williams dissimilarity update formula (Murtagh and Legendre 2011; Murtagh and Contreras 2012). In this work, the *hclust* R (<http://www.R-project.org>) function is used, which performs

hierarchical clustering and proposes a set of dissimilarity measures.

- *K-means clustering* (McQueen 1967) is a vector quantization method, originally from signal processing, that is popular for cluster analysis in machine learning. K-means clustering aims at partitioning the  $n$  observations into  $k$  clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells. The most common algorithms use an iterative refinement and alternates between assignment steps and update steps (McKay 2003). They can only be applied on intensities here because means can not be properly defined for binary variables. The *kmeans* R function is used, allowing the joint use of the Hartigan–Wong, Lloyd and McQueen algorithms (Hartigan and Wong 1979; Lloyd 1957, 1982).

Of course, both algorithms need the choice of an appropriate similarity measure, or distance, to quantify the neighborhood relation between two objects.

Clustering on binary positions vectors is first considered here. Specific similarity measures such as Ochiai, Jaccard, Dice or Russel–Rao are needed to capture the binary specificity (Plasse et al. 2007). To avoid redundancy, only two of them are used in this paper: the Jaccard and Ochiai ones. Given  $p$  binary (0=absent; 1=present) attributes, like the positions in this paper, these similarity measures between any two objects  $X$  and  $Y$  of a library are built from a general contingency table, counting the number of common attributes. Let us call  $a$  the number of common attributes when  $X = Y = 1$ ,  $b$  when  $X = 1$  and  $Y = 0$ ,  $c$  when  $X = 0$  and  $Y = 1$  and finally  $d$  when  $X = Y = 0$ . On this basis, the Jaccard and Ochiai similarity measures have both a range of 0 to 1 and are defined as follows:

$$\text{Jaccard}(X, Y) = \frac{a}{a + b + c} = \frac{|X \cap Y|}{|X \cup Y|} \quad (1)$$

$$\text{Ochiai}(X, Y) = \sqrt{\frac{a}{a + b}} \sqrt{\frac{a}{a + c}} = \frac{|X \cap Y|}{\sqrt{|X| \cdot |Y|}}. \quad (2)$$

Obviously, two spectra with many common peaks are supposed to provide high measures. The Jaccard measure can be interpreted as the size of the intersection divided by the size of the union of the sample sets; and the Ochiai measure as a geometric mean of the probabilities that if one object has the attribute, the other object has it too. Because product increases weaker than sum when only one of the terms grows, Ochiai will be really high only if both of the two proportions (probabilities) are high. It implies that the objects must share a great part of their attributes to be considered similar by Ochiai. Notice that the Ochiai

coefficient can be considered as being identical to the cosine similarity index (Holliday et al. 2002).

Clustering algorithms on recoded intensities are also implemented and use the classic euclidean distance as similarity measure. In Cartesian coordinates, if  $X = (x_1, x_2, \dots, x_p)$  and  $Y = (y_1, y_2, \dots, y_p)$  are two objects in  $\mathbb{R}^p$ , then the euclidean distance from  $X$  to  $Y$ , or from  $Y$  to  $X$ , is given by:

$$d(X, Y) = d(Y, X) = \sqrt{\sum_{i=1}^p (y_i - x_i)^2} = \|Y - X\|_2. \quad (3)$$

In this paper, the number of clusters is set to  $K = 4$  in order to correspond to the initial number of mixtures (first design) and to the initial number of blood donors (second design). Multiple combinations between Ward/K-Means algorithms, between intensities/positions vectors, between one/two/three decimals for the resolution of the GPL matrices (and between Ochiai/Jaccard measures for the binary part of the clustering) were experimented; results are discussed in Sect. 4.

### 3.5 Numerical indexes to evaluate the quality of the clustering results

To allow an objective and numerical evaluation of the quality of the clustering results, four indexes have been used: Dunn, Davies–Bouldin, Rand and Adjusted Rand indexes. The two first evaluate the homogeneity of clusters regardless of the initial groups content; the two last measure the correctness of the non-supervised clustering process according to these initial groups.

The *Dunn index* (DI) (Dunn 1973) corresponds to the ratio between the smallest distance between observations not in the same cluster and the largest intra-cluster distance. Let  $\mathbf{C} = (C_1, \dots, C_K)$  be a particular clustering partition of  $n$  objects into  $K$  disjoint clusters. Let  $\Delta_m$  be the maximum distance between observations in the cluster  $C_m$ . The Dunn index is computed as:

$$DI_{\mathbf{C}} = \min_{C_k, C_l \in \mathbf{C}, C_k \neq C_l} \left\{ \frac{\min_{i \in C_k, j \in C_l} \text{dist}(i, j)}{\max_{C_m \in \mathbf{C}} \Delta_m} \right\}. \quad (4)$$

In this work, the evaluation of inter and intra-cluster distances is based on euclidean, Ochiai or Jaccard metrics according to the nature of the data.

The *Davies–Bouldin index* (DBI) (Davies and Bouldin 1979) is an internal evaluation scheme, where the validation of how well the clustering has been done is evaluated using quantities and features inherent to the dataset. Let  $C_i$  be a cluster. Let  $x_j$  be a  $p$  dimensional feature vector of the objects assigned to cluster  $C_i$ ,  $A_i$  the medoid associated with  $C_i$  and  $|C_i|$  the cardinality of  $C_i$ . Medoids are

preferred compared with centroids in order to be able to work on binary positions. With these preliminary definitions, let  $\Gamma_i$  be the euclidean measure of variation within this cluster:

$$\Gamma_i = \sqrt{\frac{1}{|C_i|} \sum_{j=1}^{|C_i|} \|x_j - A_i\|_2^2}.$$

The euclidean example is detailed here but note that many other distance metrics can be used. Again, in the case of binary variables, the euclidean distance is replaced by the Ochiai and Jaccard ones. Let also  $\gamma(C_i, C_j)$  be a measure of separation between cluster  $C_i$  and cluster  $C_j$ , and  $a_{i,l}$  the  $l^{\text{th}}$  element of  $A_i$ , ( $l = 1, \dots, L$ ). One can write:

$$\gamma(C_i, C_j) = d(A_i, A_j) = \|A_i - A_j\|_2 = \sqrt{\sum_{l=1}^L |a_{i,l} - a_{j,l}|^2}.$$

On this basis, the Davies–Bouldin index is defined as follows:

$$DBIC \equiv \frac{1}{K} \sum_{i=1}^K \max_{j: i \neq j} \left\{ \frac{\Gamma_i + \Gamma_j}{\gamma(C_i, C_j)} \right\}. \quad (5)$$

Finally, the *Rand index* (RI) and the *Adjusted Rand index* (ARI) (Hubert and Arabie 1985) are measures of the similarity between two data clusterings, used to evaluate the quality of the classification. From a mathematical point of view, the Rand index is related to the accuracy, but is applicable even when class labels are not directly used. Given again a set of  $n$  objects  $X = \{x_1, \dots, x_n\}$  and two partitions of  $X$  to compare,  $C^1 = \{C_1^1, \dots, C_{K_1}^1\}$ , a partition into  $K_1$  subsets, and  $C^2 = \{C_1^2, \dots, C_{K_2}^2\}$ , a partition into  $K_2$  subsets, let us define the following notations:

- $m_a$  as the number of pairs of elements in  $X$  that are in the same set in  $C^1$  and in the same set in  $C^2$ ,
- $m_b$  as the number of pairs of elements in  $X$  that are in the same set in  $C^1$  and in different sets in  $C^2$ ,
- $m_c$  as the number of pairs of elements in  $X$  that are in different sets in  $C^1$  and in the same set in  $C^2$ ,
- $m_d$  as the number of pairs of elements in  $X$  that are in different sets in  $C^1$  and in different sets in  $C^2$ .

The Rand index is defined as:

$$RI = \frac{m_a + m_d}{m_a + m_b + m_c + m_d} \quad (6)$$

Intuitively,  $m_a + m_d$  can be considered as the number of agreements between  $C^1$  and  $C^2$ , and  $m_b + m_c$  as the number of disagreements. The Rand index has a value between 0 and 1, with 0 indicating that the two data clusters do not

agree on any pair of points and 1 indicating that the data clusters are exactly the same.

The Adjusted Rand index (ARI) is the corrected-for-chance version of the Rand index. Though the Rand Index may only yield a value between 0 and 1, the Adjusted Rand Index can yield negative values if the index is less than the expected index. Following the same notations  $m_a$ ,  $m_b$ ,  $m_c$  and  $m_d$ , it is defined as (Santos and Embrechts 2009):

$$\text{ARI} = \frac{\binom{n}{2}(m_a + m_d) - [(m_a + m_b)(m_a + m_c) + (m_b + m_d)(m_c + m_d)]}{\binom{n}{2}^2 - [(m_a + m_b)(m_a + m_c) + (m_b + m_d)(m_c + m_d)]} \quad (7)$$

Note that DI, RI and ARI have to be as high as possible to represent good clustering partitions. DBI has to be as small as possible.

## 4 Results and discussion

In this section, clustering methods and related quality indexes are used to assess the performance of different versions of COSY spectra. The first upcoming subsection proposes numerical outputs based on the indexes defined in Sect. 3.5 for both designs and for combinations of experiments described in Sects. 3.3 and 3.4. Then, in Sect. 4.2, comparisons between clustering results based on 2D COSY spectra and on corresponding  $^1\text{H-NMR}$  spectra are shown. This section will allow to demonstrate on the two designs that, indeed, COSY spectra include additive relevant information for spectral classification as compared to  $^1\text{H-NMR}$  ones (providing consequently a better MIC). By doing this, a statistically-proved response to the initial key question (does supplementary information include relevant information?) is provided.

### 4.1 Results for COSY spectra

The numerical results are available in Table 2 for the first experimental design and for the more complex second one. Note that only a part of all the considered combinations are shown to avoid too much redundancy and to improve clarity.

Globally, the results are very promising with many RI (and ARI) greater than 0.7, particularly with bucketed data when one or two decimals in the GPL matrices are used. Unlabeled individuals are well grouped together, with no prior knowledge, and this according to the presence of numerous potential noise factors (sampling, time replicates, degradation, changes in temperature, etc...).

Considering first the position-based partitions, particular groups are already well isolated in the majority of cases

(for instance, the third cell culture mixture in the first design was known to be significantly different from the others). An example is given in Fig. 2a.

With the recoded intensities-based partitions, the four mixtures are generally well recovered by the algorithms in spite of the sampling procedure and time repetitions. For the first design, the best clustering result is obtained with the one-decimal GPL matrix, thus underlying the impor-

tance of the “bucketing” step: Fig. 2b shows that there is only one error during the blind clustering process and RI and ARI indexes are maximized (RI=0.927 and ARI=0.853 in Table 2). The conclusion is the same for the second design concerning the blood donors (RI = 0.932 and ARI = 0.804 in Table 2).

More generally, Dunn and Davies–Bouldin indexes are subject to significative variations between experiments and are not very satisfactory for some combinations (several DI less than 0.5 and DBI greater than 1.5 in Table 2). This means that obtained clusters are not always compact and well separated. However, to judge if the clusters conform to the reality, i.e. if members in a cluster correspond to members of initial groups, Rand and Adjusted Rand indexes prevail because they are built on the degree of agreement and disagreement between groups (here between the reality and each of the clustering partitions).

In conclusion, and based on the two experimental designs, the clustering results show that COSY spectra appear to be statistically robust and contain informative signal which helps to succeed in distinguishing groups of different spectra. In other words, COSY spectra are enough informative so that the signal connected to the initial groups is distinguished from the noise. Working on positions gives satisfactory results but not better ones than working on intensities. Moreover, these results show that the 2D “bucketing” step is important and probably necessary to solve misalignment problems (as it is the case in 1D). This additional information can be profitable for further robust statistical analysis, as biomarker discovery.

### 4.2 Comparisons with corresponding $^1\text{H-NMR}$ spectra

Besides the convincing performances of COSY spectra, this section intends to demonstrate that the additional information contained in 2D spectra (compared to 1D) is relevant and crucial by improving the quality of the



**Table 2** Numerical clustering performances according to different versions of COSY spectra

Signal	Dec.	Algorithm	Distance	<i>T</i>	DI	DBI	RI	ARI
<i>First design</i>								
Positions	1	Ward	Jaccard	909	0.712	1.466	<b>0.914</b>	<b>0.811</b>
Positions	1	Ward	Ochiai	909	0.712	1.466	<b>0.914</b>	<b>0.811</b>
Positions	2	Ward	Jaccard	2348	<b>0.857</b>	1.688	0.757	0.420
Positions	2	Ward	Ochiai	2348	<b>0.827</b>	1.688	0.757	0.420
Positions	3	Ward	Jaccard	3250	<b>0.893</b>	1.722	0.601	0.189
Positions	3	Ward	Ochiai	3250	<b>0.892</b>	1.722	0.601	0.189
Intensities	1	Ward	Euclidean	909	0.298	<b>0.541</b>	<b>0.927</b>	<b>0.853</b>
Intensities	1	K-means	Euclidean	909	0.298	<b>0.541</b>	<b>0.927</b>	<b>0.853</b>
Intensities	2	Ward	Euclidean	2348	0.239	1.249	0.669	0.609
Intensities	2	K-means	Euclidean	2348	0.311	1.356	0.657	0.501
Intensities	3	Ward	Euclidean	3250	0.439	1.540	0.588	0.425
Intensities	3	K-means	Euclidean	3250	0.434	1.614	0.597	0.439
<i>Second design</i>								
Positions	1	Ward	Jaccard	1106	0.796	1.569	<b>0.937</b>	<b>0.825</b>
Positions	1	Ward	Ochiai	1106	0.722	1.569	<b>0.937</b>	<b>0.825</b>
Positions	2	Ward	Jaccard	4172	<b>0.945</b>	1.800	0.772	0.401
Positions	2	Ward	Ochiai	4172	<b>0.899</b>	1.800	0.772	0.401
Positions	3	Ward	Jaccard	6686	<b>0.981</b>	1.792	0.650	0.171
Positions	3	Ward	Ochiai	6686	<b>0.963</b>	1.792	0.650	0.171
Intensities	1	Ward	Euclidean	1106	0.419	<b>0.643</b>	<b>0.932</b>	<b>0.804</b>
Intensities	1	K-means	Euclidean	1106	0.419	<b>0.643</b>	<b>0.932</b>	<b>0.804</b>
Intensities	2	Ward	Euclidean	4172	0.689	1.592	0.789	0.422
Intensities	2	K-means	Euclidean	4172	0.706	1.742	0.735	0.288
Intensities	3	Ward	Euclidean	6686	0.778	1.668	0.578	0.055
Intensities	3	K-means	Euclidean	6686	0.765	1.886	0.647	0.135

Bold indexes are highlighted the best performances for each index

The column “Dec.” denotes the number of decimal(s) for the data in the GPL matrix; *T* number of observations in the corresponding GPL, *DI* Dunn index, *DBI* Davies–Bouldin index, *RI* Rand index, *ARI* adjusted Rand index

clustering results. It is always difficult to directly compare objects of different dimensions. To be able to compare <sup>1</sup>H-NMR and COSY spectra, some pre-processing steps are necessary to upgrade the one-dimensional spectra before performing clustering.

#### 4.2.1 Upgrade of <sup>1</sup>H-NMR spectra

Some very powerful and complete tools are available to pre-process, transform and interpret <sup>1</sup>H-NMR data, for example in Xia and Wishart (2010) and Vanwinsberghe (2005). Here, the goal is not to focus on sophisticated pre-processing methods for <sup>1</sup>H-NMR data, but just to upgrade them in order to obtain comparable pre-processing levels for both 1D and COSY data. To do so, one-dimensional spectra were subject to the following steps after phase and baseline correction:

- elimination of negative intensities,
- deletion of the water spectral zone between 4.5 and 5.5 ppm,

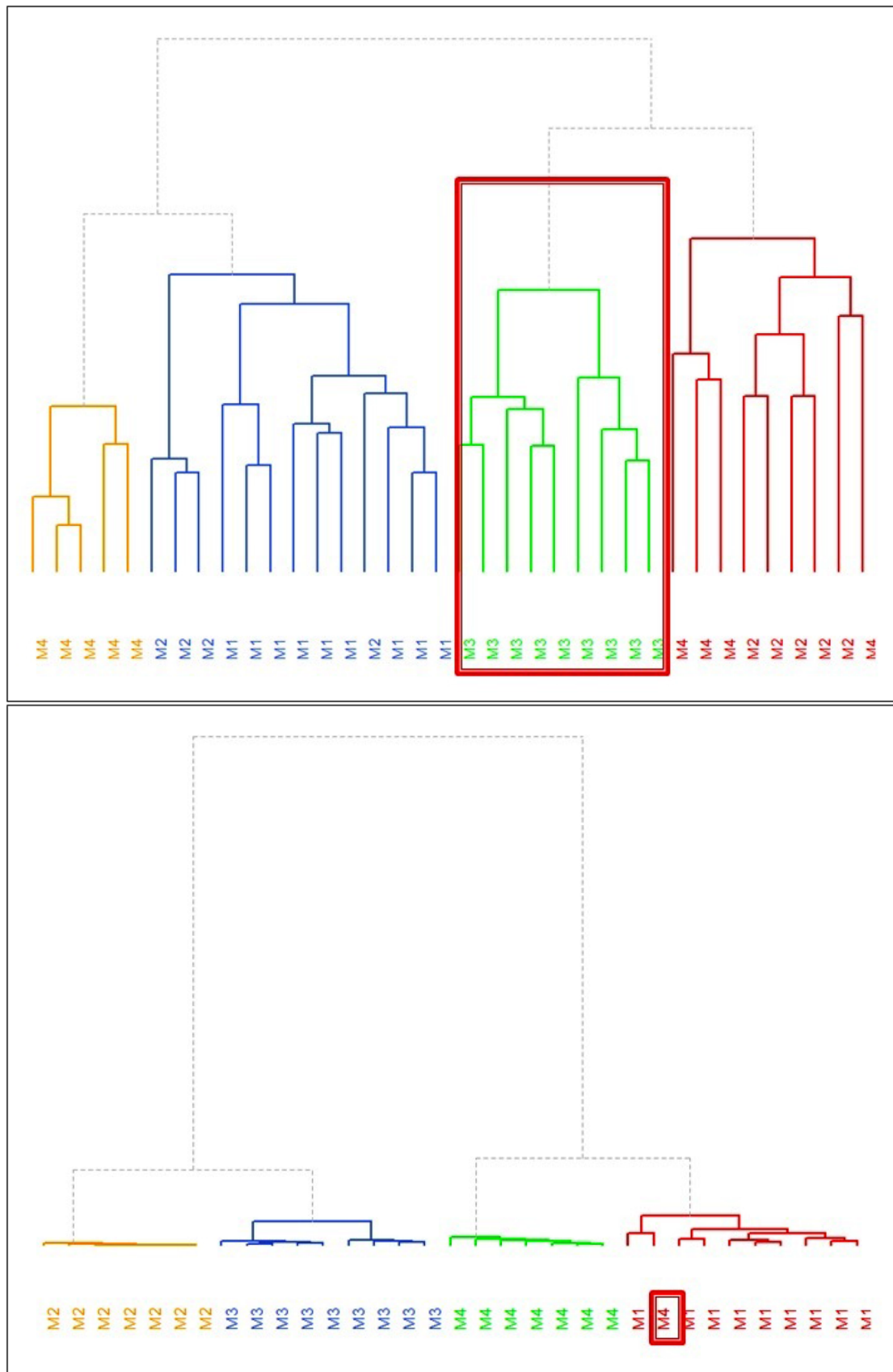
- application of the constant sum transformation,
- “bucketing” by using the same number of decimals (for the ppm coordinates) than for the COSY experiments,
- outlier deletion.

In a general way, it is quite intuitive that 2D spectra require less pre-processing to obtain an acceptable visualization of the signal and of potential biomarkers. Finally, note that the reasoning based on the positions makes no sense when working with 1D spectra. Consequently, clustering processes were only performed on 1D normalized intensities (raw intensities, “bucketed” with two decimals and with only one decimal).

#### 4.2.2 Comparisons

The clustering results obtained with <sup>1</sup>H-NMR spectra are available in Table 3 for the two designs.

If one compares these results with Rand and Adjusted Rand indexes in Table 2, it appears clearly that these



**Fig. 2** **a** Illustrative dendrogram of the clustering on positions (first design, Ward algorithm with Ochiai distance, dec = 2, before outlier deletion). M1 to M4 denote the four different initial mixtures. Different colours are attributed to different clusters; **b** Illustrative dendrogram of the clustering on intensities (first design, Ward algorithm, dec = 1, after outlier deletion)

indexes are mainly higher when the clustering algorithms are performed on COSY spectra. Resulting partitions of the spectra are better for COSY: they are more in agreement with the real initial groups. Thus, this proves, for the two experiments, that the additional spectral dimension does provide relevant information to improve the clustering results and, consequently, to improve the amount of captured signal and the spectral informative content (MIC).

## 5 Conclusion and further works

In this article, an advanced clustering approach is proposed to evaluate and quantify the “Metabolomic Informative Content” (as defined in the introduction) of NMR spectral data in metabolomics. It is applied on both  $^1\text{H}$ -NMR and COSY spectra, using both qualitative input (binary positions, linked with the absence or presence of a peak or metabolite) or quantitative inputs (intensities). More precisely, the choice of this clustering approach is to be related to the presence of initial groups in the signal. In the two real experimental designs detailed in Sect. 2.2, four different cell culture mixtures and blood samples from four different donors were respectively involved. The elements of these groups were “mixed” via several noisy factors:

sampling, time replicates, delays, deterioration, risk of bacterial contamination, etc...The goal of the clustering processes was to blindly recover these initial groups using all the individual unlabeled spectra. And the final quality of these processes informs us about the quantity of captured signal and can be finally viewed as a measure of MIC.

This work shows that COSY spectra allow to well discriminate groups linked with different signals, proving that they are successful in spite of the noisy factors. It also shows that the clustering processes perform better with COSY spectra than with  $^1\text{H}$ -NMR ones, thus confirming a higher MIC for the two-dimensional spectra.

More precisely, this work first shows some pre-processing steps for the data analyst in order to handle 2D COSY data, including the construction of the GPL matrix and an ad hoc 2D bucketing step which allows to control the data resolution. It then demonstrates that COSY spectra provide very satisfying clustering results: in other words, in spite of the multiple noise factors, the informative part of the signal can be well discovered and, finally, the final clusters are mainly in concordance with the initial groups. The clustering methodology also highlights that 2D bucketing, by aggregating data according to reduced numbers of decimals for the ppm coordinates, helps to improve the amount of captured information when using COSY data.

Furthermore, this paper also provides a comparison between COSY and  $^1\text{H}$ -NMR spectra, based on the quality of respective clustering results. It is demonstrated that COSY spectra provide more MIC about the groups than corresponding upgraded  $^1\text{H}$ -NMR spectra, thus proving the importance and the relevance of the additional dimension

**Table 3** Numerical clustering performances for 1D data

Signal	Dec.	Algorithm	Distance	$T$	DI	DBI	RI	ARI
<i>First design</i>								
Intensities	1	Ward	Euclidean	201	<b>0.927</b>	<b>0.120</b>	<b>0.908</b>	<b>0.807</b>
Intensities	1	K-means	Euclidean	201	0.560	<b>0.463</b>	<b>0.835</b>	0.667
Intensities	2	Ward	Euclidean	2001	0.712	<b>0.426</b>	0.717	0.516
Intensities	2	K-means	Euclidean	2001	0.601	0.739	0.732	0.476
Raw intensities		Ward	Euclidean	199967	0.362	1.048	0.606	0.544
Raw intensities		K-means	Euclidean	199967	0.284	1.182	0.410	0.196
<i>Second design</i>								
Intensities	1	Ward	Euclidean	207	<b>0.981</b>	<b>0.688</b>	0.704	0.233
Intensities	1	K-means	Euclidean	207	0.334	<b>0.663</b>	0.702	0.230
Intensities	2	Ward	Euclidean	2054	<b>0.901</b>	0.986	0.704	0.233
Intensities	2	K-means	Euclidean	2054	0.635	1.119	<b>0.740</b>	<b>0.290</b>
Raw intensities		Ward	Euclidean	205034	0.670	1.074	0.588	0.017
Raw intensities		K-means	Euclidean	205034	0.447	1.321	0.675	0.145

Bold indexes are highlighted the best performances for each index

The column “Dec.” denotes the number of decimal(s) for the data in the GPL matrix;  $T$  number of observations in the corresponding GPL,  $DI$  Dunn Index,  $DBI$  Davies–Bouldin Index,  $RI$  Rand Index,  $ARI$  Adjusted Rand Index

of COSY to go further in NMR-based metabolomics studies. Note that this methodology can be generalized for comparing any spectrometric tools when groups are present and of interest in the experimental design.

These promising results have to be confirmed on faster versions of COSY experiments, in order to deal with more competitive acquisition times. The methodology can also be applied to heteronuclear two-dimensional spectroscopy tools, like HSQC spectra. Ultimately, once the MIC of a given 2D tool or technology is verified, the research of discriminating zones or biomarkers, which represents the main objective of most metabolomics studies, can benefit from the advantageous use of 2D-NMR spectra instead of the more traditional  $^1\text{H}$ -NMR ones (for instance, in PLS-DA, OPLS-DA or more advanced machine learning techniques).

**Softwares** The raw data were firstly processed with the Bruker Topspin 2.1 software. Peak lists were extracted using ACD/Labs 12.00 (ACD/NMR processor). For manipulating the 1D and COSY data (pre-processings steps), generating the GPL matrices and performing the clustering processes, the R software environment was exclusively used (<http://www.R-project.org>).

**Acknowledgments** This work was supported by the FNRS from which P. de Tullio is senior research associate. Support from the IAP Research Network P7/06 of the Belgian State (Belgian Science Policy) is also gratefully acknowledged.

#### Compliance with ethical standards

**Conflict of Interest** Authors declare that they have no conflict of interest.

**Human and Animal Rights and Informed Consent** This study analyzes collected data which involved human participants who had provided informed consent.

## References

- Akitt J.W., Mann B.E. (2000). NMR and Chemistry (Manual), Cheltenham UK, Stanley Thornes. p. 287.
- Aue, W. P., Bartholdi, E., & Ernst, R. R. (1976). Two-dimensional spectroscopy. Application to nuclear magnetic resonance. *The Journal of Chemical Physics*, 64, 2229–2246.
- Bruschweiler, R., & Bingol, K. (2011). Deconvolution of chemical mixtures with high complexity by NMR consensus trace clustering. *Analytical Chemistry*, 83(19), 7412–7417.
- Bruschweiler, R., Bingol, K., Bruschweiler-Li, L., & Li, D.-W. (2014). Customized metabolomics database for the analysis of NMR  $^1\text{H}$ - $^1\text{H}$  TOCSY and  $^{13}\text{C}$ - $^1\text{H}$  HSQC-TOCSY Spectra of Complex Mixtures. *Analytical Chemistry*, 86(11), 5494–5501.
- Craig, A., Cloarec, O., Holmes, E., Nicholson, J. K., & Lindon, J. C. (2006). Scaling and normalization effects in NMR spectroscopic metabolomic data sets. *Analytical Chemistry*, 78, 2262–2267.
- Davies, D. L., & Bouldin, D. W. (1979). A cluster separation measure. *IEEE Transactions on Pattern Analysis and Machine Intelligence, PAMI-1*(2), 224–227.
- Dunn, J. C. (1973). A fuzzy relative of the ISODATA process and its use in detecting compact well-separated clusters. *Journal of Cybernetics*, 3(3), 32–57.
- Giraudeau, P., Remaud, G., & Akoka, S. (2009). Evaluation of Ultrafast 2D NMR for quantitative analysis. *Analytical Chemistry*, 81(1), 479–484.
- Hartigan, J. A., & Wong, M. A. (1979). A K-means clustering algorithm. *Applied Statistics*, 28, 100–108.
- Holliday, J. D., Hu, C. Y., & Willett, P. (2002). Grouping of coefficients for the calculation of inter-molecular similarity and dissimilarity using 2D fragment bit-strings. *Combinatorial Chemistry and High Throughput Screening*, 5(2), 155–166.
- Hubert, L., & Arabie, P. (1985). Comparing partitions. *Journal of Classification*, 2(1), 193–218.
- Iman, R. L. (2008). *Latin hypercube sampling*. New York: Wiley.
- Keeler, J. (2010). *Understanding NMR Spectroscopy* (2nd ed., pp. 190–191). New York: Wiley.
- Le Guennec, A., Giraudeau, P., & Caldarelli, S. (2014). Evaluation of fast 2D NMR for metabolomics. *Analytical chemistry*, 86(12), 5946–5954.
- Lloyd S. P., Least squares quantization in PCM, Technical Note, Bell Laboratories, IEEE Transactions on Information Theory 28, pp. 128–137 (1957, 1982).
- MacKay, D. (2003). *An Example Inference Task: Clustering, Information Theory, Inference and Learning Algorithms* (pp. 284–292). Cambridge: Cambridge University Press.
- MacQueen J. B. (1967). Some Methods for classification and Analysis of Multivariate Observations. In *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability*, (vol 1), University of California Press, (pp. 281–297) .
- Mao, X., & Ye, C. (1997). Phase-shift presaturation for water peak suppression in biomolecular NMR experiments, Science in China. *Series C, Life sciences*, 40(4), 345–350.
- Marion, D., & Bax, A. (1988). Baseline distortion in real-fourier-transform NMR spectra. *Journal of Magnetic Resonance* (1969), 79(2), 252–356.
- Murtagh F., Legendre P., Ward's hierarchical clustering method: clustering criterion and agglomerative algorithm, arXiv preprint [arXiv:1111.6285](https://arxiv.org/abs/1111.6285) (2011).
- Murtagh, F., & Contreras, P. (2012). Algorithms for hierarchical clustering: an overview. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2(1), 86–97.
- Nicholson, J., Connelly, J., Lindon, J. C., & Holmes, E. (2002). Metabonomics: a generic platform for the study of drug toxicity and gene function. *Nature Reviews Drug Discovery*, 1, 153–161.
- Plasse, M., Niang, N., Saporta, G., Villeminot, A., & Leblond, L. (2007). Combined use of association rules mining and clustering methods to find relevant links between binary rare attributes in a large data set. *Computational Statistics and Data Analysis*, 52(1), 596–613.
- Queiroz Junior, L. H. K., Ferreira, A. G., & Giraudeau, P. (2013). Optimization and practical implementation of ultrafast 2D NMR experiments. *Quimica Nova*, 36(4), 577–581.
- Rasmussen, L. G., Savorani, F., Larsen, T. M., Dragsted, L. O., Astrup, A., & Engelsen, S. B. (2011). Standardization of factors that influence human urine metabolomics. *Metabolomics*, 7(1), 71–83.
- Rousseau R., Statistical contribution to the analysis of metabolomic data in  $^1\text{H}$ -NMR spectroscopy, PhD Thesis, UCL, <http://hdl.handle.net/2078.1/75532> (2011).
- Santos, J. M., & Embrechts, M. (2009). *On the use of the adjusted rand index as a metric for evaluating supervised classification, Artificial Neural Networks, ICANN 2009*. Berlin: Springer.
- Sousa, S. A., Magalhaes, A., & Castro Ferreira, M. M. (2013). Optimized bucketing for NMR spectra. *Chemometrics and Intelligent Laboratory Systems*, 122, 93–102.
- Vanwinsberghe J., Bubble: development of a matlab tool for automated  $^1\text{H}$ -NMR data processing in metabolomics, Master's thesis, Université de Strasbourg (2005).

- Vega-Vazquez, M., Cobas, J. C., & Martin-Pastor, M. (2010). Fast multidimensional localized parallel NMR spectroscopy for the analysis of samples. *Magnetic Resonance in Chemistry*, 48(10), 749–752.
- Ward, J. H. (1963). Hierarchical Grouping to optimize an objective function. *Journal of American Statistical Association*, 58(301), 236–244.
- Xi, Y., deRopp, J. S., Viant, M., Woodruff, D., & Yu, P. (2007). Automated screening for metabolites in complex mixtures using 2D COSY NMR spectroscopy. *Metabolomics*, 2(4), 221–233.
- Xia, J., & Wishart, D. (2010). MetPA: a web-based metabolomics tool for pathway analysis and visualization. *Bioinformatics*, 26(18), 2342–2344.
- Yun, K., Sunghyouk, P., Jongheon, S., & Dong-Chan, O. (2013). Application of  $^{13}\text{C}$ -labeling and  $^{13}\text{C}$ - $^{13}\text{C}$  COSY NMR experiments in the structure determination of a microbial natural product. *Archive of Pharmacal Research*,. doi:[10.1007/s12272-013-0254-8](https://doi.org/10.1007/s12272-013-0254-8).