



Universiteit  
Leiden  
The Netherlands

## **Metabolomics of biofluids : from analytical tools to data interpretation**

Nevedomskaya, E.

### **Citation**

Nevedomskaya, E. (2011, November 23). *Metabolomics of biofluids : from analytical tools to data interpretation*. Retrieved from <https://hdl.handle.net/1887/18135>

Version: Corrected Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/18135>

**Note:** To cite this publication please use the final published version (if applicable).

# Chapter

Cross-platform analysis of  
longitudinal data in metabolomics

*Nevedomskaya E., Mayboroda O.A., Deelder A.M.*

Molecular Biosystems, 2011, DOI: 10.1039/c1mb05280b

6

## ABSTRACT

Metabolic profiling is considered to be a very promising tool for diagnostic purposes, for assessing nutritional status and response to drugs. However, it is also evident that human metabolic profiles have a complex nature, influenced by many external factors. This, together with the understanding of the difficulty to assign people to distinct groups and a general move in clinical science towards personalized medicine, raises the interest to explore individual and variable metabolic features for each individual separately in longitudinal study design. In the current paper we have analyzed a set of metabolic profiles of a selection of six urine samples per person from a set of healthy individuals by  $^1\text{H}$  NMR and reversed-phase UPLC-MS. We have demonstrated that the method for recovery of individual metabolic phenotypes can give complementary information to another established method for analysis of longitudinal data—multilevel component analysis. We also show that individual metabolic signatures can be found not only in  $^1\text{H}$  NMR data, as has been demonstrated before, but also even more strongly in LC-MS data.

## INTRODUCTION

Metabolomics is a post-genomics technology, the aim of which is “profiling metabolism in complex systems”.(1) The reasoning behind metabolomics experiments is that metabolites, compared to genes, transcripts and proteins, offer the closest representation of the phenotype.(2) Thus, they can contain valuable information on a disease development, in contrast to genetics that gives insight into predisposition to a disease. Conducting this type of research assumes the existence of specific biomarkers or metabolic signatures that can distinguish between pathological states. The potential of metabolomics to reveal signatures of pathological conditions has been demonstrated on a number of neurological disorders (Huntington’s disease,(3) Parkinson’s disease,(4) multiple sclerosis(5) *etc.*), various cancers (ovarian and breast,(6) pancreatic,(7) colorectal(8) and others), cardiological abnormalities (*e.g.* ischemia(9)) and many other diseases.

It is, however, evident, that metabolic profiles reflect not only the disease/healthy state of the organism but, as a representation of a given phenotype, they strongly depend on factors such as gender, age and our daily habits, like, for example, diet, drugs and alcohol intake.(10,11) Another extremely important factor affecting metabolic profiles is gut microbiota. It has been shown that even for genetically identical laboratory animals gut microflora is influenced by environment and dietary factors.(12–14) For humans, in which genetic variation is enormous, gut microflora is much more diverse and is additionally affected by factors such as, for instance, stress.(15)

Therefore, human metabolic profiles are highly dependent on environmental factors and may vary from day to day due to turbulent conditions of our fast and highly stressful modern lives. A way to overcome possible negative effects of this variability on the interpretation of metabolomics data is multiple sampling over time per individual. The main advantage of this approach is the possibility to get an insight into the biological processes, which are usually missed by a simple, static comparison of “diseased vs. healthy”. The feasibility of the longitudinal sampling has been demonstrated more than once in toxicology experiments in animals,(16) in human intervention studies(17) and in monitoring cyclic dynamics in healthy women.(18)

A number of statistical methods that can deal with longitudinal(19) or paired(20) data exist. One of the methods suitable for the analysis of such data is the multilevel approach that separates levels of variation present in the data into inter- and intra-individual.(21) Moreover, it has been shown that, despite the multiple sources of variability, present in human biofluids and particularly in urine, there are constant individual metabolic signatures, which are probably to a great extent determined by genetics.(22) Assfalg et al. used a combination of established statistical methods for individual classification, or in

other words person recognition, on the basis of Nuclear Magnetic Resonance (NMR) spectra of urine. The core of this approach was a variant of Principal Component Analysis (PCA), which was an innovative tool for face recognition 20 years ago,(23) and proved to be innovative for recognition of human urine metabolic signatures today. The existence of stable personal metabolic phenotypes is linked to the idea of homeostasis or, more precisely, to the idea that individual, self-regulatory, genetically controlled mechanisms maintain the homeostasis at any price. Consequently, the disruption of homeostasis means the beginning of a disease development.(24;25) Thus, monitoring of individual metabolic signatures, which represent dynamic, time-correlated changes of the phenotype, might ultimately develop into a preferred approach for practical personalized medicine.

Thus, there are different statistical methods that can be applied to the metabolic profiles of multiple samples per individual and allow having various perspectives on the data. Besides, an enormous chemical diversity of metabolites has resulted in a broad spectrum of analytical approaches used in metabolomics, especially with regards to mass spectrometry. To select an auxiliary to NMR MS method, one has to make a choice between gas chromatography (GC), capillary electrophoresis (CE) and liquid chromatography (LC). Choosing the latter, one still remains with a range of possibilities like, for example, reverse phase or hydrophilic interaction liquid chromatography (rpLC or HILIC).

In the current research we wanted to demonstrate that multilevel component analysis (MCA) and person recognition can be used in parallel, and that the information retrieved by the two methods is complimentary and together they can form a toolbox for analysis of longitudinal datasets. To this end we analyzed a set of urine samples from 8 healthy individuals (each of them donated 6 samples) by  $^1\text{H}$  NMR and reversed phase UPLC-MS, which requires relatively simple sample preparation and is one of the most widely used MS techniques in metabolomics. MCA and person recognition methods were applied to the data. We here show that individual metabolic phenotypes can be identified not only on the basis of  $^1\text{H}$  NMR spectra, as has been shown before, but also on the basis of LC-MS data. We also demonstrate the information extracted from this type of designed study using the two statistical approaches, based on diverse sources of variation in the data. The difference in information content of the data from the two analytical techniques is analyzed and discussed as well.

## MATERIALS AND METHODS

**Samples.** Urine samples were collected, after written consent, from 8 self-declared healthy individuals from the same working environment (Leiden University Medical Center, the Netherlands). The volunteers were equally divided into men and women, aged

between 25 and 45 years old, all Caucasian. Each volunteer provided 6 urine samples of the first morning urine after over-night fasting on 6 different days (5 consecutive weekdays and one after the weekend). No diet restrictions were implied; none of the subjects was taking medication. Samples were collected in sterile 15 ml polypropylene tubes, kept at 4 °C, frozen within 8 h of collection, and stored for approximately 2 weeks at -20 °C until the measurement. In total 48 urine samples were analyzed.

**Sample preparation.** Frozen samples were thawed at room temperature and vortexed before use.

*Sample preparation for <sup>1</sup>H NMR experiments.* Aliquots of urine sample (1000 μl) were centrifuged at 3000g for 15 min at 4 °C to remove any precipitate. 600 μl of each sample were transferred to a 96 deep-well plate, further preparation was automated using the Bruker Sample Track system and a Gilson 215 robotic system. Here 540 μl urine were added to 60 μl of pH 7.0 sodium phosphate buffer (0.2 M) in 10% D<sub>2</sub>O containing 0.53 mM sodium 3-trimethylsilyl-tetradeuteriopropionate (TSP) and 0.26 mM NaN<sub>3</sub>, thoroughly mixed and transferred to a new 96 deep-well plate. Samples were centrifuged at 3000g for 5 min to remove any solid debris. A modified Gilson 215 robot was used to transfer 565 μl of sample from the plate into 5 mm SampleJet NMR tubes.

*Sample preparation for LC-MS experiments.* 150 μl of each urine sample were mixed with 450 μl of water and subsequently centrifuged at 3000 rpm for 10 min. 5 μl of sample was used for injection.

#### **Data acquisitions.**

*<sup>1</sup>H NMR experiments.* All <sup>1</sup>H NMR experiments were performed on a 600 MHz Bruker Avance II spectrometer (Bruker BioSpin, Karlsruhe, Germany) equipped with a 5 mm TCI cryogenic probe head with Z-gradient system and automatic tuning and matching. Temperature calibration was done prior to the measurements using the method of Findeisen *et al.*(40)

One-dimensional <sup>1</sup>H NMR spectra were recorded at 300 K using the first increment of a NOESY(41) pulse sequence with presaturation ( $\gamma$ B1 = 50 Hz) during a relaxation delay of 4 s and a mixing time of 10 ms for efficient water suppression. A total of 32 768 data points were recorded with 32 scans covering a sweep width of 12 336 Hz. The free induction decay (FID) was zero-filled to 65 536 complex data points prior to Fourier transformation and an exponential window function was applied with a line broadening factor of 1.0 Hz.

A sample <sup>1</sup>H NMR spectrum can be found in Supplementary Materials, Figure S1a.

**LC-MS experiments.** The samples were analyzed on a UPLC-ESI-UHR-ToF system. The injection scheme was randomized and included quality control samples (mix of all of the urine samples, prepared in the same way as the individual samples), as well as a set of

analytical standards (mix of pesticides, see Supplementary Materials, Table S1) to ensure the robustness of the workflow and to evaluate the analytical variability. Quality control (QC) and analytical standards were injected at the beginning and at the end of the sequence, as well as every four biological samples. In total 28 QC runs were acquired. The UPLC (Ultimate 3000 RS tandem LC system, Dionex, Amsterdam, The Netherlands) was equipped with a pre-column (Acclaim 120 C18, 5 mm, 120 Å, 2 × 10 mm) and two analytical columns (Acclaim RSLC 120 C18, 2.2 mm, 120 Å, 2.1 × 100 mm) working alternatively to speed up the acquisition series. The UPLC flow was set at 400 µl min<sup>-1</sup> and the mobile phases were water + 0.1% formic acid v/v (Phase A) and methanol+0.1% formic acid v/v (Phase B). The gradient was as follows: 1 min 0% phase B, then in 1 min to 10% phase B, held for 1 min at 10% phase B, and subsequently in 6,5 min to 100% phase B and held for 3 min at 100% phase B. Before each chromatographic run, a calibrant solution of sodium formate was injected in Flow Injection Analysis mode.

The ESI-UHR-ToF (maXis, Bruker Daltonics, Bremen, Germany) was operated in the positive ionization mode and acquired data in the mass range from *m/z* 50 to 1500 with a spectra rate of 1 Hz. The capillary was set at 2500 V, the End Plate offset at -500 V, the Nebulizer gas at 2 bar and the dry gas at 8 l min<sup>-1</sup> at 180 °C.

A sample LC-MS chromatogram can be found in Supplementary Materials, Figure S1b.

**Data pre-processing.** *<sup>1</sup>H NMR data pre-processing.* All spectra were manually phase- and baseline-corrected using Topspin 2.1 (Bruker BioSpin, Karlsruhe, Germany) and automatically referenced to TSP signal (0.0 ppm). Each spectrum was integrated (binned) using 0.0095 ppm integral regions between 0.5 and 10 ppm, the residual water and urea region between 4.5 and 6 ppm was excluded, resulting in 842 bin regions used for the analysis. To account for any difference in concentration between the samples, each spectrum was normalized to its total area and subsequently by Probabilistic Quotient Normalization (PQN) (42) using average spectrum as a reference.

**LC-MS data pre-processing.** All data files were recalibrated on the masses of sodium formate clusters. The alignment of chromatograms and peak picking was performed using open-source XCMS software (The Scripps Research Institute, La Jolla, CA).(43) Finding peaks was performed using the “centWave” algorithm with *m/z* deviation set to 5 ppm, and the scan range between 20 and 700 scans. Grouping of peaks was done with parameters minsamp set to 28 (number of QC samples) and bandwidth to 10. Retention time correction was done with default parameters. The resulting table included the detected ion features and their peak areas. The peaks were filtered on the basis of QC samples: the peak was retained in the analysis if it was present in all the QC samples and relative standard deviation of the area in QC samples was less than 20%. The final table contained 965 ion

features, which areas were normalized on total areas of the samples and subsequently by PQN(42) with an average of QC samples taken as a reference.

The consistency of the data and the absence of column-bias were checked using Principal Component Analysis (Supplementary Materials, Figure S2).

**Statistical data analysis.** *Principal Component Analysis* was performed on logarithmically transformed, mean-centered and unit variance scaled data using the NIPALS algorithm.(44)

*Person recognition.* The person recognition approach used in the current paper was based on the previously published classification method.(22) Among the classification methods used by Assfalg *et al.* the combination of Principal Component Analysis (PCA) for data reduction and canonical discriminant analysis (CA) was chosen as the most effective one. The accuracy of classification was assessed using test-set validation: in each round of validation one randomly selected sample per donor was taken out into the test-set, and a model was built based on the remaining samples. The test-set samples were projected into the PCA–CA subspace and classified according to the minimum distance to the mean of the discriminated groups. The resulting class labels were compared to the real ones and the number of correct classifications was evaluated. The validation was performed in 1000 rounds and the results averaged throughout all the rounds (Supplementary Materials, Figure S3a).

Recognition accuracies were also assessed in 100 rounds of Subject ID permutations and compared with the actual accuracies, statistical significance of the difference was assessed using the Mann–Whitney test.

*Multilevel Components Analysis (MCA).* MCA is an effective method for separating the variation between- and within- individuals and analyzing them by different submodels. The method was implemented as described by Jansen *et al.*(21) PCA were performed on the data matrix corresponding to the between-individual variation and on the within-individual variation for each individual separately (Supplementary Materials, Figure S3b).

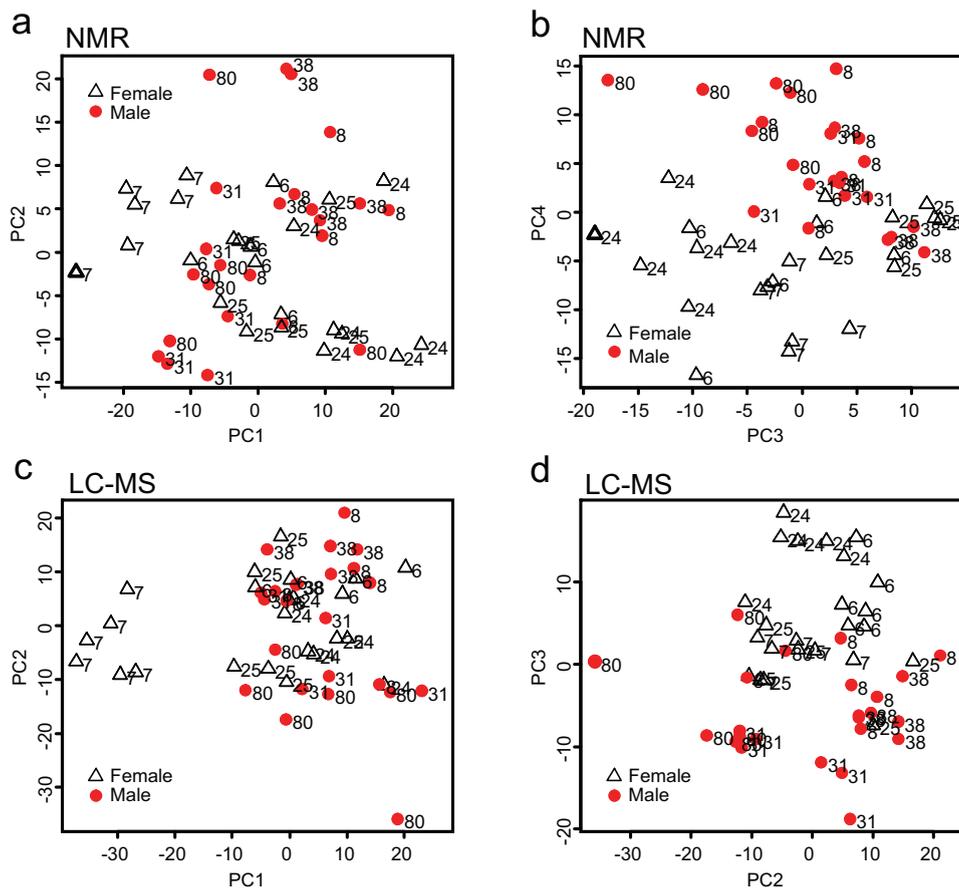
All the data-preprocessing and statistical analysis were performed in a R statistical software environment (<http://www.r-project.org/>) using in-house developed scripts and open-source packages.

## RESULTS

A set of longitudinal urine samples from healthy individuals was analyzed by <sup>1</sup>H NMR and rpUPLC-MS and subsequently by various statistical methods in order to compare the information that can be retrieved from the data by different analytical and statistical approaches. PCA was performed on <sup>1</sup>H NMR and LC-MS data. The scores plot of the first

two principal components for each of the techniques revealed some clustering by individuals; however no separation between the genders was observed (Fig. 1a and c). Difference between the genders was visible on the scores plot of the third and fourth principal components in the case of the  $^1\text{H}$  NMR data (Figure 1b), and of the second and third principal components in the case of the LC-MS data (Figure 1d). At a first glance, there appeared to be a similarity between the score plots of the first two principal components on  $^1\text{H}$  NMR and LC-MS data, for example, subject 7 is separated from the rest of the people. A way to give a numerical value to this similarity is the use of the RV-coefficient, which is a multivariate extension of correlation coefficient. This has already been used before for estimation of the overlap of metabolomics data matrices, but in that case both matrices were derived from MS-based experiments.(26) For all the 8 principal components of the PCA the RV-coefficient was found to be not that high, not exceeding 0.46. The RV-coefficient does not increase anymore after the fourth component, which explains 44 and 53% of the variation in the  $^1\text{H}$  NMR and LC-MS data respectively (Table 1). Hence, the relative positions of data points in the PCA subspace are different for  $^1\text{H}$  NMR and LC-MS, with little overlap.

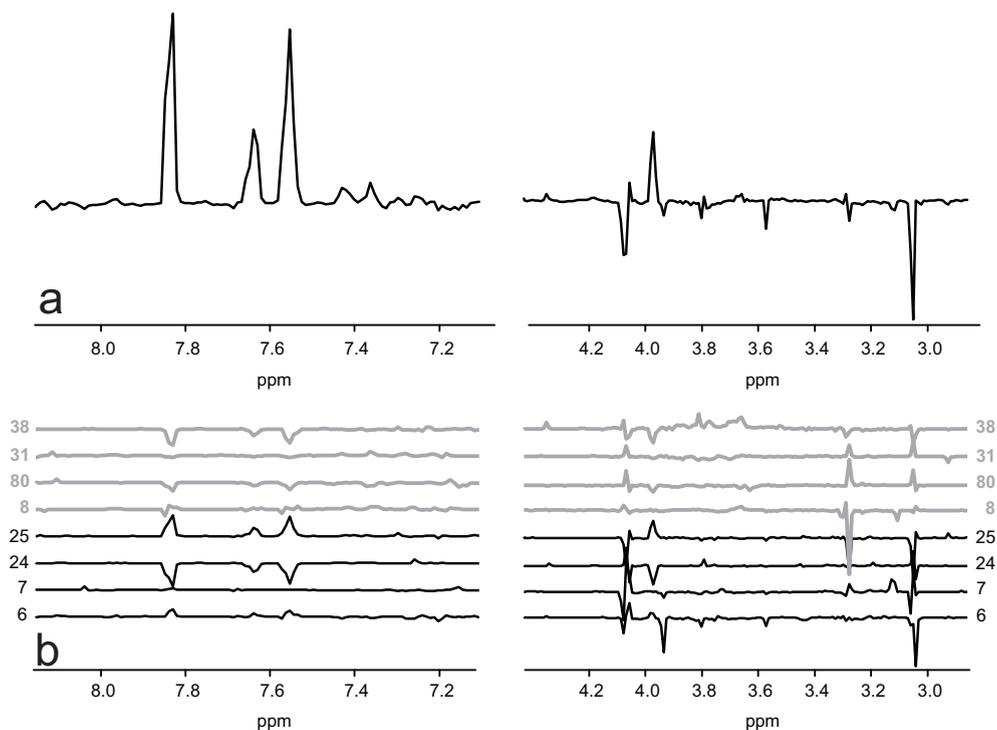
It was evident from the PCA analysis that there are different sources of variation present in the data, which this method is not capable of separating. One of the ways to dissect variations present at different levels in the data (*e.g.* between and within individuals) is to use multilevel analysis which has been successfully applied for a number of applications in social sciences, geography, public health (27) and recently also in metabolomics.(20,21) This method was applied to the data so that the variation at two levels—between individuals and within each of the individuals—was explored. It allows identification of spectral regions or peaks variable between individuals and peaks/regions variable between the time points for each individual. In Figure 2 the loading plots for the between and within-individual models are shown on the example of the  $^1\text{H}$  NMR data, demonstrating those variable areas. The RV-coefficient, calculated for the results of multilevel analysis (on the between-individual score matrices), is higher than that for the PCA analysis, but is very stable even with the growth of the explained variance, again indicating that the two analytical techniques explain different relations between the samples (Table 2).



**Figure 1.** Scores plots of the PCA analysis of  $^1\text{H}$  NMR and LC-MS data from urine samples of 8 individuals sampled at six different time points. Samples are labeled by individuals' IDs. Triangles represent urine samples of females, dots of males. (a) First two principal components of the PCA on NMR data cover 16.1 and 9.6% of variation respectively. (b) Third and fourth principal components of the PCA on NMR data cover 7.3% and 6.6% of variation, respectively. (c) First two principal components of the PCA on LC-MS data cover 21.1 and 12.3% of variation, respectively. (d) Second and third principal components of the PCA on LC-MS data, the third component covers 7.5% of variation.

**Table 1. Summary of principal component analysis performed on NMR and LC-MS data and multivariate correlation (RV-coefficient), calculated on the resulting score matrices.**

PC No.	Explained variation, %		RV
	NMR	LC-MS	
1	16.1	21.2	0.2
2	25.7	33.5	0.24
3	33	41	0.28
4	39.6	47.9	0.4
5	44.3	53	0.41
6	48.3	57.7	0.42
7	52.1	61.9	0.44
8	55.6	65.9	0.46



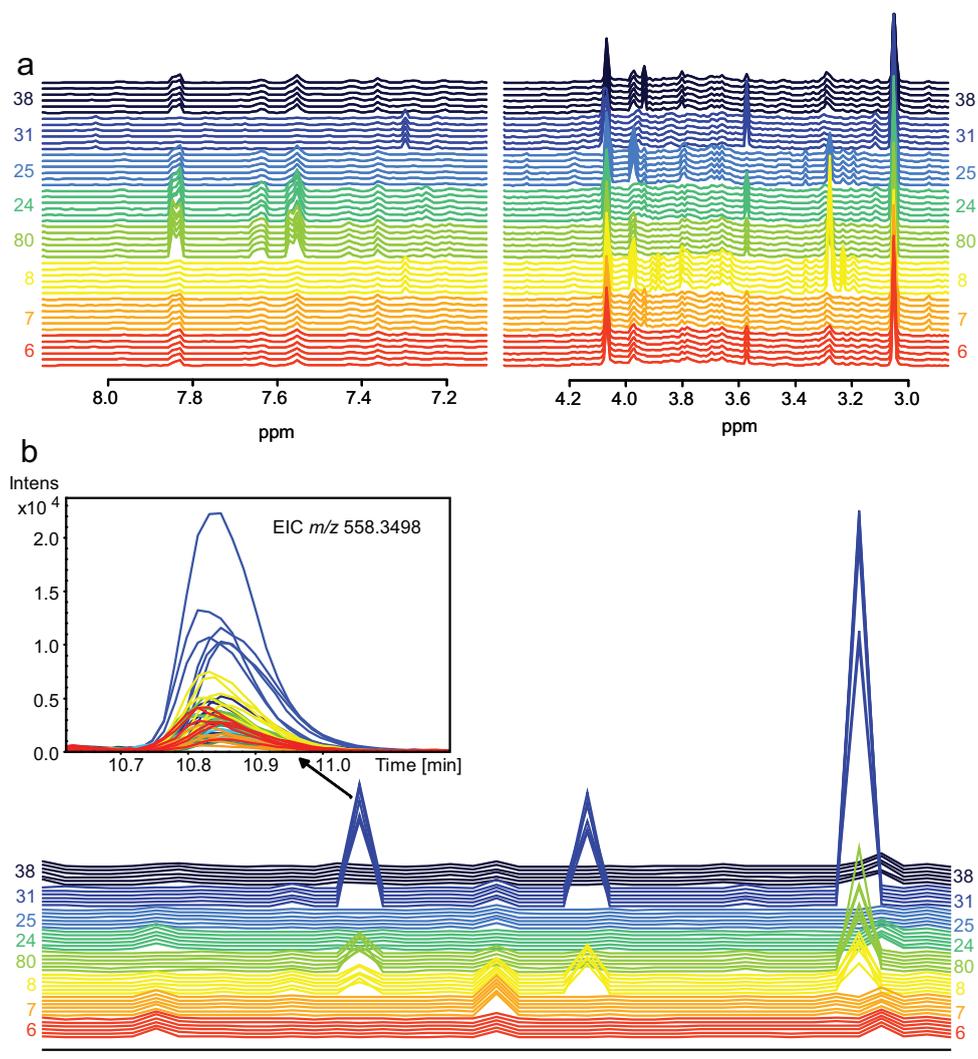
**Figure 2. Loadings plot of the multilevel component analysis of  $^1\text{H}$  NMR data of urine samples from 8 individuals sampled at six different time points: (a) loadings of the first component of the between-individual model, (b) loadings of the first component of the within-individual model for each person, colored by gender (grey—male, black—female) and labeled by individuals' IDs.**

Depending on the research question, one might not be interested in the most variable regions for each individual, but in the most constant ones, person-specific features. Those unique patterns can be used to recognize each person from the rest. The existence of such fingerprints in  $^1\text{H}$  NMR data and a method to assess them were demonstrated previously (22) using an innovative combination of classical statistical methods—PCA and canonical discriminant analysis with thorough validation. We performed person recognition on the  $^1\text{H}$  NMR data evaluating the accuracy with which each of the people is predicted. The recognition accuracy ranged from 59.5 to 99.5% which matches the estimated probability of correct classification for the same number of samples in the model described in the previous work.(22) The mean recognition was 84%, which is quite high taking into account that in each validation step the model is built only on 5 spectra. The accuracy of recognition was also calculated on the set with permuted person labels and it appeared to be significantly lower (mean accuracy 13%,  $p$ -value < 0.001) than the real recognition results (Supplementary Materials, Figure S4a).

One of the advantages of the person recognition method is that it is possible to perform back-projection of scores in the canonical subspace into the PCA scores subspace and then into the original variables. As a result of this procedure individual metabolic phenotypes are obtained (Figure 3a). These metabolic phenotypes represent the characteristic spectral regions for each person and are, unlike the original profiles, easily clustered by *e.g.* hierarchical clustering per person (Supplementary Materials, Figure S5).

**Table 2. Summary of multilevel component analysis (between individuals) performed on NMR and LC-MS data and multivariate correlation (RV-coefficient), calculated on the resulting score matrices.**

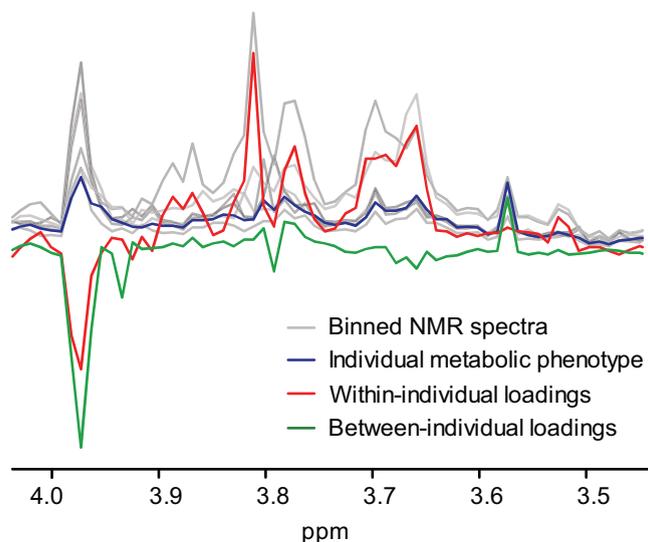
PC No.	Explained variation, %		
	NMR	LC-MS	RV
1	39.4	49.8	0.64
2	70	78.1	0.6
3	87.8	88.1	0.6
4	93.8	94.6	0.62
5	96.9	97.3	0.63
6	99	99.2	0.65



**Figure 3. Individual metabolic phenotypes found within urine samples of 8 individuals sampled at six different time points. (a) Based on the  $^1\text{H}$  NMR spectra. (b) Based on LC-MS data; one of the differential peaks ( $m/z$  of 558.3498) was used to construct the corresponding extracted ion chromatograms (inlet). Colored and labeled by individuals' IDs.**

To illustrate the results and differences of person recognition and multilevel analysis methods we here show an example for one of the participants, using the  $^1\text{H}$  NMR spectra and its analysis. In Figure 4 the original binned  $^1\text{H}$  NMR spectra are shown together with

the individual metabolic profile of the person and his within- and between-individual loadings. As can be seen from the picture, there are peaks that contribute differently to the different types of analyses. For instance, the peak at 3.97 ppm appears in all of the six urine spectra, but its intensity is variable, so it contributes to both the individual metabolic profile and within-individual loadings. This peak is found to be also characteristic for separation between individuals, contributing to the between-individual loadings. Another peak at 3.81 ppm appears only in one of the six spectra, and thus it contributes to the within-individual loadings, indicating its variability through the time course of the study, however, it is not characteristic of the person (does not appear in the individual metabolic profile), neither is responsible for the separation between the subjects. Another example is the peak at 3.57 ppm which is present in all the samples and has quite a consistent intensity; it does not contribute to the within individual loadings, but is characteristic for the individual and drives the separation between the individuals.



**Figure 4.** Comparison of the statistical analyses of the urinary  $^1\text{H}$  NMR data from one individual. Grey lines show original (binned)  $^1\text{H}$  NMR spectra. Blue line indicates the mean individual metabolic profile, red is within-individual loadings of the first component, green—between-individual loadings of the first component.

The individual metabolic profiles were found to exist not only in  $^1\text{H}$  NMR spectra, but also in UPLC-MS data. Moreover, the accuracies of the person recognition performed on LC-MS data were found to be significantly higher than those derived from  $^1\text{H}$  NMR data: the range of accuracies was between 92 and 100% and the mean recognition was 98.3% (Supplementary Materials, Figure S4b). Recognition accuracies in the randomized

experiment were significantly lower (mean 12.6%,  $p$ -value < 0.001) than the real values. In the same way as for  $^1\text{H}$  NMR, the LC-MS peaks specific for an individual can be found by back-projection (Figure 3b) and traced back in the original data. Indeed they show a differential profile across individuals. The recognition accuracies based on  $^1\text{H}$  NMR and LC-MS are not correlated across the subjects (correlation coefficient is -0.012): the individuals that are relatively badly recognized on the basis of  $^1\text{H}$  NMR may be well-recognized on the basis of LC-MS and *vice versa*.

Thus, in the current study a small longitudinal cohort of urine samples from healthy individuals was analyzed, and even with this limited sample set it was evident that various sources of variation are confounded. There are statistical methods available to extract this variation separately and examine the data from a different perspective. Those methods (multilevel component analysis and person recognition) can be applied to both  $^1\text{H}$  NMR and LC-MS data. The use of multiple analytical platforms also widens the information extracted from a study.

## DISCUSSION

The importance of a personalized approach in health-related research is being widely accepted by the scientific community, (28) and metabolic profiling is recognized as a valuable tool for personalized medicine.(29) However in the majority of metabolomics studies a traditional “case-control” design is applied, although it is well-known that the definition of these groups, and especially the group of healthy individuals, is very vague.(30) In contrast, the definition of the physical individuality is very clear and monitoring an individual reaction to perturbations and its development in time is a promising approach for medicine and pharmacology.

In metabolomics the advantages of the longitudinal study design that implies multiple sampling per individual have clearly been demonstrated for interventional studies where the dynamic response of the organism to the drug or other substance can be monitored.(31;32) It has also been demonstrated that “classical” data analysis methods used in metabolomics, such as PCA and PLS-DA, are suboptimal(20) for such dynamic data and other methods, separating levels of variation, should be used. The longitudinal design offers possibilities for differential analysis; depending on the question of interest one may focus on differences between the subjects, on variations within each subject or identification of unique profiles for each subjects.

To illustrate some of the possibilities for the analysis of longitudinal metabolomics data we applied a series of statistical methods to the  $^1\text{H}$  NMR and LC-MS data of urines of 8 individuals, 4 females and 4 males, each contributing 6 urine samples for the study. As no

special diet or life style restrictions were applied, it was obvious that the data would contain a lot of variation due to day-to-day differences, between-subjects diversity and certain grouping of the samples due to, for example, gender, age *etc.* This, indeed, was confirmed by PCA, which summarizes the variation present in the data. The clustering according to person was evident; however day-to-day variation for most of the people was much higher than the differences between individuals. Gender distinction was also present, but not in the first two components of PCA, suggesting that the between-gender difference is overruled by all the other sources of variation.

In the PCA analysis LC-MS data showed more variation covered in the first principal components, than the  $^1\text{H}$  NMR data (Table 1). There was some similarity visible for the position of the data points along the first principal components in the two datasets; however RV-coefficients calculated on the principal component subspaces for the two techniques were rather low (Table 1). This suggests that the two analytical methods reveal different biological phenomena reflected in the metabolic composition of the same biological samples.

PCA modeling has shown that the metabolic data with underlying design contain information from different sources—from the variation between the subjects, as well as the variation between the samples for each person. There is a class of multilevel statistical models which can separate the data into levels and as such are perfectly applicable to our data. In the case of MCA, applied in the current study, the overall variation present in the study was divided into between-individual and within-individual and separate analysis was performed on each block. This method reveals spectral regions differential between the people, as well as regions which are variable for each of the people through the time course of the analysis.

Another method used, namely, person recognition, focuses on different parts of the spectra, which are consistent for the individual and thus characteristic. Before, this method was successfully applied to  $^1\text{H}$  NMR spectra, revealing the existence of individual metabolic signatures, which were found to be extremely stable over time and could be largely explained by genetics.(22,33) We successfully applied the described method to our data and observed a recognition accuracy corresponding to the number of samples used. We also demonstrated on an example how the individual metabolic profiles give information complementary to that derived from multilevel analysis.

As can be seen from the analysis, different levels of variation can be recovered from a set of spectra. The choice of an appropriate method for statistical analysis should be based on the question posed by the investigation. The right answers can only be derived from a carefully designed and analyzed study. Consequently, longitudinal design offers possibilities

for real personalized medicine such as exploration of the effects not averaged between people, elimination of day-to-day variation, focusing on intrinsic individual properties reflected in the metabolic composition of urine.

The person recognition strategy, to the best of our knowledge, has so far not been applied to LC-MS data, which is one of the most commonly used analytical techniques in metabolomics experiments.(34) One study evaluating the amount of “personalized” information present in a set of LC-MS data has been conducted,(35) however this study only described the features, unique for an individual (i.e. appearing in one set of the spectra), but not the unique patterns of the features as in a person recognition approach. Thus our report appeared to be the first ever attempt to do the person recognition analysis on LC-MS data. Despite the fact that the LC-MS has somewhat lower analytical reproducibility than  $^1\text{H}$  NMR, (36) person recognition accuracy was substantially higher in the case of LC-MS (all individuals were recognized with accuracy more than 92%, compared to 59% in  $^1\text{H}$  NMR). There was absolutely no correlation between recognition of people in  $^1\text{H}$  NMR and in LC-MS again pointing at the fact that the two techniques most probably provide different information concerning the samples.

Of course, the differences in the metabolome coverage between  $^1\text{H}$  NMR and LC-MS come as no surprise.  $^1\text{H}$  NMR is a universal approach capable of detecting all the compounds that contain hydrogens, whereas MS-based methods are more targeted due to the selectivity of the separation and ionization techniques used.(36) On the other hand,  $^1\text{H}$  NMR has slightly lower sensitivity in comparison to MS. Thus, there is a certain “bias” in metabolomics experiments performed on a single analytical platform: the coverage of the metabolites is either limited by the sensitivity, or by the separation method. Thus, the observed “personalized” content of LC-MS data in comparison to  $^1\text{H}$  NMR might be a result of such analytical bias. This, however, has to be further explored; here we can make only a few assumptions about what is driving this difference.

In general,  $^1\text{H}$  NMR-based metabolomics studies result in a systemic view on the studied phenomenon, due to the fact that a lot of the detected compounds are related to energy metabolism: TCA cycle intermediates, amino acids, *etc.* These molecules are highly abundant in biofluids, are also day-to-day variable depending on the diet and are also involved in many biochemical pathways. The latter means that they change in many states of the organism, which leads to the problem of “usual suspects”(37) with many of the same metabolites discovered to be differential in a number of conditions.(38)

In contrast to NMR, LC-MS is more specific due to the inherent selectivity of the separation method and high sensitivity of the detection. Reversed-phase UPLC-MS is highly suitable for separation of medium polar and non-polar compounds.(39) Most of the

molecules related to energy metabolism, amino and other organic acids will not be retained. Thus, the part of the metabolome, covered by rpUPLC-MS, might be less affected by diet and gut microflora and might provide a closer approximation of the phenotype.

This phenomenon certainly would need more extensive investigation and might be an extremely important issue in the decision how to conduct a study using a certain analytical platform, depending on the study design and the question of interest.

In total, the amount of biologically relevant information that can be derived from metabolomics experiments is enormous. However, the quality of this information depends on a clear definition of the goals and the study design as well as on the selection of the analytical platform and the subsequent statistical analysis. All of these factors are extremely important for obtaining successful results and generating a relevant hypothesis. Consequently, performing costly, labor-intensive metabolomics experiments with the sole aim to distinguish “diseased from healthy” might be seen as a suboptimal use of manpower, instrumental resources and, the most importantly patient material. The power of a longitudinal design and the flexibility of various statistical methods to analyze such a design may open new possibilities. Individual metabolic signatures, that represent dynamic, time-correlated changes of phenotype, may actually be used as a phenotype-readout essential for practical personalized medicine.

## CONCLUSIONS

In the metabolomics-related literature somewhat controversial ideas are present: on the one hand that metabolites can provide unique diagnostic information, and on the other hand that their concentrations are very sensitive to non-systemic external factors and vary even from day to day for one individual. However, it has been demonstrated that highly individual metabolic signatures exist in for instance urine, on top of which the other variation is superimposed. The available methods for the analysis of time-resolved data can focus either on variation between people or within the time course for an individual. In the current paper we have demonstrated the use and complementarity of the extracted information of some of these statistical methods on a set of data from healthy individuals.

We have also shown that the detection of individual metabolic profiles is not solely the property of  $^1\text{H}$  NMR, but is also possible based on UPLC-MS data, interestingly—even with a higher accuracy. Based on this limited data set, it would appear that the parallel analysis of  $^1\text{H}$  NMR and LC-MS indicates that the two techniques explain different phenomena in the data. The higher accuracy of person recognition in LC-MS further suggests that the method might be more sensitive to unique, individual-specific features, while  $^1\text{H}$  NMR might reflect a more systemic response.

## ACKNOWLEDGEMENTS

The authors would like to thank all the volunteers for the supplied samples, Sibel Göröler M.Sc. and Ing. Bart Schoenmaker for the analytical work, Dr Hartmut Schäfer for his help with implementing the person recognition method, Dr Paul J. Hensbergen for fruitful discussion.

## REFERENCES

1. Lindon, J. C.; Holmes, E.; Bollard, M. E.; Stanley, E. G.; Nicholson, J. K. Metabonomics technologies and their applications in physiological monitoring, drug safety assessment and disease diagnosis. *Biomarkers* 2004, 9 (1), 1-31.
2. Nicholson, J. K.; Connelly, J.; Lindon, J. C.; Holmes, E. Metabonomics: a platform for studying drug toxicity and gene function. *Nat. Rev. Drug Discov.* 2002, 1 (2), 153-161.
3. Underwood, B. R.; Broadhurst, D.; Dunn, W. B.; Ellis, D. I.; Michell, A. W.; Vacher, C.; Mosedale, D. E.; Kell, D. B.; Barker, R. A.; Grainger, D. J.; Rubinsztein, D. C. Huntington disease patients and transgenic mice have similar pro-catabolic serum metabolite profiles. *Brain* 2006, 129 (Pt 4), 877-886.
4. Bogdanov, M.; Matson, W. R.; Wang, L.; Matson, T.; Saunders-Pullman, R.; Bressman, S. S.; Flint, B. M. Metabolomic profiling to develop blood biomarkers for Parkinson's disease. *Brain* 2008, 131 (Pt 2), 389-396.
5. Simone, I. L.; Federico, F.; Trojano, M.; Tortorella, C.; Liguori, M.; Giannini, P.; Picciola, E.; Natile, G.; Livrea, P. High resolution proton MR spectroscopy of cerebrospinal fluid in MS patients. Comparison with biochemical changes in demyelinating plaques. *J. Neurol. Sci.* 1996, 144 (1-2), 182-190.
6. Slupsky, C. M.; Steed, H.; Wells, T. H.; Dabbs, K.; Schepansky, A.; Capstick, V.; Fought, W.; Sawyer, M. B. Urine metabolite analysis offers potential early diagnosis of ovarian and breast cancers. *Clin. Cancer Res.* 2010, 16 (23), 5835-5841.
7. Nishiumi, S.; Shinohara, M.; Ikeda, A.; Yoshie, T.; Hatano, N.; Kakuyama, S.; Mizuno, S.; Sanuki, T.; Kutsumi, H.; Fukusaki, E.; Azuma, T.; Takenawa, T.; Yoshida, M. Serum metabolomics as a novel diagnostic approach for pancreatic cancer. *Metabolomics* 2010, 6 (4), 518-528.
8. Wang, H.; Tso, V. K.; Slupsky, C. M.; Fedorak, R. N. Metabolomics and detection of colorectal cancer in humans: a systematic review. *Future. Oncol.* 2010, 6 (9), 1395-1406.
9. Barba, I.; de Leon, G.; Martin, E.; Cuevas, A.; Aguade, S.; Candell-Riera, J.; Barrabes, J. A.; Garcia-Dorado, D. Nuclear magnetic resonance-based metabolomics predicts exercise-induced ischemia in patients with suspected coronary artery disease. *Magn Reson. Med.* 2008, 60 (1), 27-32.
10. Holmes, E.; Loo, R. L.; Stampler, J.; Bictash, M.; Yap, I. K.; Chan, Q.; Ebbels, T.; de, I. M.; Brown, I. J.; Veselkov, K. A.; Daviglus, M. L.; Kesteloot, H.; Ueshima, H.; Zhao, L.; Nicholson, J. K.; Elliott, P. Human metabolic phenotype diversity and its association with diet and blood pressure. *Nature* 2008, 453 (7193), 396-400.
11. Slupsky, C. M.; Rankin, K. N.; Wagner, J.; Fu, H.; Chang, D.; Weljie, A. M.; Saude, E. J.; Lix, B.; Adamko, D. J.; Shah, S.; Greiner, R.; Sykes, B. D.; Marrie, T. J. Investigations of the effects of gender, diurnal variation, and age in human urinary metabolomic profiles. *Anal Chem* 2007, 79 (18), 6995-7004.
12. Friswell, M. K.; Gika, H.; Stratford, I. J.; Theodoridis, G.; Telfer, B.; Wilson, I. D.; Mcbain, A.

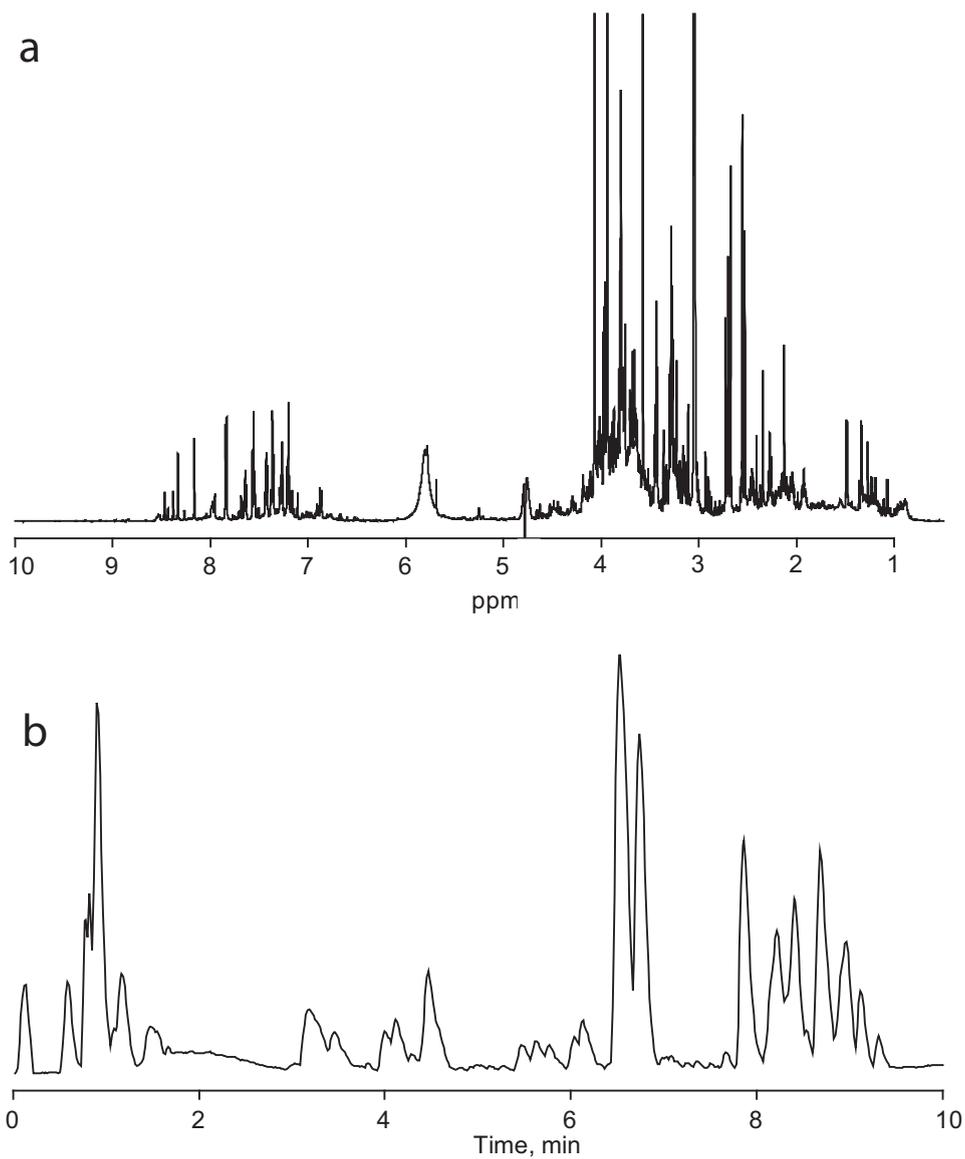
- J. Site and Strain-Specific Variation in Gut Microbiota Profiles and Metabolism in Experimental Mice. *Plos One* 2010, 5 (1).
13. Martin, F. P.; Sprenger, N.; Montoliu, I.; Rezzi, S.; Kochhar, S.; Nicholson, J. K. Dietary modulation of gut functional ecology studied by fecal metabolomics. *J. Proteome Res.* 2010, 9 (10), 5284-5295.
  14. Phipps, A. N.; Stewart, J.; Wright, B.; Wilson, I. D. Effect of diet on the urinary excretion of hippuric acid and other dietary-derived aromatics in rat. A complex interaction between diet, gut microflora and substrate specificity. *Xenobiotica* 1998, 28 (5), 527-537.
  15. Rezzi, S.; Martin, F. P.; Alonso, C.; Guilarte, M.; Vicario, M.; Ramos, L.; Martinez, C.; Lobo, B.; Saperas, E.; Malagelada, J. R.; Santos, J.; Kochhar, S. Metabotyping of Biofluids Reveals Stress-Based Differences in Gut Permeability in Healthy Individuals. *Journal of Proteome Research* 2009, 8 (10), 4799-4809.
  16. Keun, H. C.; Ebbels, T. M. D.; Bollard, M. E.; Beckonert, O.; Antti, H.; Holmes, E.; Lindon, J. C.; Nicholson, J. K. Geometric trajectory analysis of metabolic responses to toxicity can define treatment specific profiles. *Chemical Research in Toxicology* 2004, 17 (5), 579-587.
  17. van Velzen, E. J.; Westerhuis, J. A.; van Duynhoven, J. P.; van Dorsten, F. A.; Grun, C. H.; Jacobs, D. M.; Duchateau, G. S.; Vis, D. J.; Smilde, A. K. Phenotyping tea consumers by nutrikinetic analysis of polyphenolic end-metabolites. *J. Proteome Res.* 2009, 8 (7), 3317-3330.
  18. Wallace, M.; Hashim, Y. Z. H. Y.; Wingfield, M.; Culliton, M.; McAuliffe, F.; Gibney, M. J.; Brennan, L. Effects of menstrual cycle phase on metabolomic profiles in premenopausal women. *Human Reproduction* 2010, 25 (4), 949-956.
  19. Smilde, A. K.; Westerhuis, J. A.; Hoefsloot, H. C. J.; Bijlsma, S.; Rubingh, C. M.; Vis, D. J.; Jellema, R. H.; Pijl, H.; Roelfsema, F.; van der Greef, J. Dynamic metabolomic data analysis: a tutorial review. *Metabolomics* 2010, 6 (1), 3-17.
  20. Westerhuis, J. A.; van Velzen, E. J. J.; Hoefsloot, H. C. J.; Smilde, A. K. Multivariate paired data analysis: multilevel PLSDA versus OPLSDA. *Metabolomics* 2010, 6 (1), 119-128.
  21. Jansen, J. J.; Hoefsloot, H. C. J.; van der Greef, J.; Timmerman, M. E.; Smilde, A. K. Multilevel component analysis of time-resolved metabolic fingerprinting data. *Analytica Chimica Acta* 2005, 530 (2), 173-183.
  22. Assfalg, M.; Bertini, I.; Colangiuli, D.; Luchinat, C.; Schafer, H.; Schutz, B.; Spraul, M. Evidence of different metabolic phenotypes in humans. *Proceedings of the National Academy of Sciences of the United States of America* 2008, 105 (5), 1420-1424.
  23. Turk, M.; Pentland, A. Eigenfaces for Recognition. *Journal of Cognitive Neuroscience* 1991, 3 (1), 71-86.
  24. van der, G. J.; Stroobant, P.; van der, H. R. The role of analytical sciences in medical systems biology. *Curr. Opin. Chem Biol.* 2004, 8 (5), 559-565.
  25. van der Greef, J.; Smilde, A. Symbiosis of chemometrics and metabolomics: past, present, and future. *JOURNAL OF CHEMOMETRICS* 2005, 19 (5-7), 376-386.
  26. Smilde, A. K.; van der Werf, M. J.; Bijlsma, S.; van der Werff-van-der Vat; Jellema, R. H. Fusion of mass spectrometry-based metabolomics data. *Analytical Chemistry* 2005, 77 (20), 6729-6736.
  27. Diez-Roux, A. V. Multilevel analysis in public health research. *Annual Review of Public Health* 2000, 21, 171-192.
  28. Weston, A. D.; Hood, L. Systems biology, proteomics, and the future of health care: toward predictive, preventative, and personalized medicine. *J. Proteome Res.* 2004, 3 (2), 179-196.

29. Schnackenberg, L. K.; Kaput, J.; Beger, R. D. *Metabolomics: a tool for personalizing medicine? Personalized Medicine* 2008, 5 (5), 495-504.
30. Elliott, R.; Pico, C.; Dommels, Y.; Wybranska, I.; Hesketh, J.; Keijer, J. *Nutrigenomic approaches for benefit-risk analysis of foods and food components: defining markers of health. Br. J. Nutr.* 2007, 98 (6), 1095-1100.
31. Holmes, E.; Antti, H. *Chemometric contributions to the evolution of metabonomics: mathematical solutions to characterising and interpreting complex biological NMR spectra. Analyst* 2002, 127 (12), 1549-1557.
32. van der Greef, J.; Hankemeier, T.; McBurney, R. N. *Metabolomics-based systems biology and personalized medicine: moving towards n=1 clinical trials? Pharmacogenomics* 2006, 7 (7), 1087-1094.
33. Bernini, P.; Bertini, I.; Luchinat, C.; Nepi, S.; Saccenti, E.; Schafer, H.; Schutz, B.; Spraul, M.; Tenori, L. *Individual Human Phenotypes in Metabolic Space and Time. J. Proteome Res.* 2009.
34. Lindon, J. C.; Nicholson, J. K. *Analytical technologies for metabonomics and metabolomics, and multi-omic information recovery. Trac-Trends in Analytical Chemistry* 2008, 27 (3), 194-204.
35. Johnson, J. M.; Yu, T. W.; Strobel, F. H.; Jones, D. P. *A practical approach to detect unique metabolic patterns for personalized medicine. Analyst* 2010, 135 (11), 2864-2870.
36. Lindon, J. C.; Nicholson, J. K. *Spectroscopic and Statistical Techniques for Information Recovery in Metabonomics and Metabolomics. Annual Review of Analytical Chemistry* 2008, 1, 45-69.
37. Robertson, D. G. *Metabonomics in toxicology: A review. Toxicological Sciences* 2005, 85 (2), 809-822.
38. Dunn, W. B.; Broadhurst, D. I.; Atherton, H. J.; Goodacre, R.; Griffin, J. L. *Systems level studies of mammalian metabolomes: the roles of mass spectrometry and nuclear magnetic resonance spectroscopy. Chem Soc. Rev.* 2011, 40 (1), 387-426.
39. Want, E. J.; Wilson, I. D.; Gika, H.; Theodoridis, G.; Plumb, R. S.; Shockcor, J.; Holmes, E.; Nicholson, J. K. *Global metabolic profiling procedures for urine using UPLC-MS. Nature Protocols* 2010, 5 (6), 1005-1018.
40. Findeisen, M.; Brand, T.; Berger, S. A <sup>1</sup>H-NMR thermometer suitable for cryoprobes. *Magn Reson. Chem* 2007, 45 (2), 175-178.
41. Kumar, A.; Ernst, R. R.; Wuthrich, K. *A two-dimensional nuclear Overhauser enhancement (2D NOE) experiment for the elucidation of complete proton-proton cross-relaxation networks in biological macromolecules. Biochem. Biophys. Res. Commun.* 1980, 95 (1), 1-6.
42. Dieterle, F.; Ross, A.; Schlotterbeck, G.; Senn, H. *Probabilistic quotient normalization as robust method to account for dilution of complex biological mixtures. Application in <sup>1</sup>H NMR metabonomics. Anal Chem* 2006, 78 (13), 4281-4290.
43. Smith, C. A.; Want, E. J.; O'Maille, G.; Abagyan, R.; Siuzdak, G. *XCMS: processing mass spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and identification. Anal Chem* 2006, 78 (3), 779-787.
44. Wold, H. *Estimation of principal components and related models by iterative least squares. In Multivariate Analysis, Krishnaiah, P. R., Ed.; Academic Press: New York, 1966; pp 391-420.*

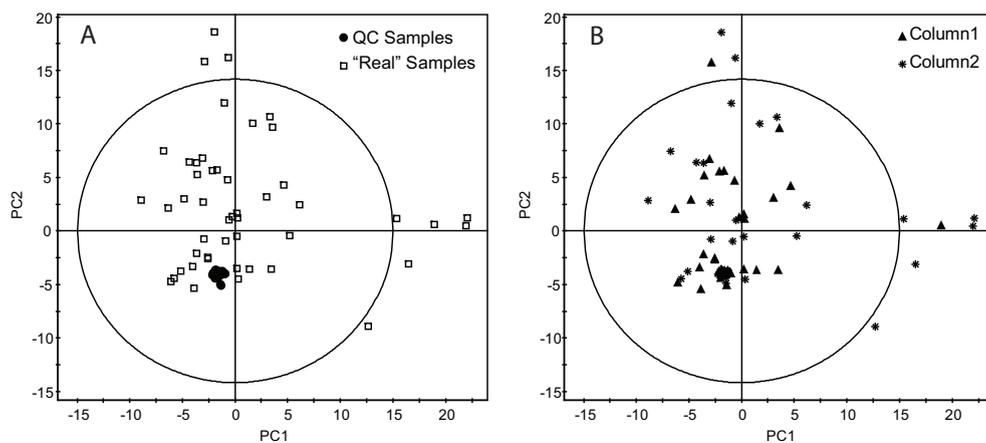
## SUPPLEMENTARY MATERIALS

**Table S1. The mix of pesticides used as the analytical standard.**

Name	Molecular formula	m/z [M+H]	Retention time, min
Pymetrozine	C <sub>10</sub> H <sub>11</sub> N <sub>5</sub> O	218.1036	4.7
Formetanate	C <sub>11</sub> H <sub>15</sub> N <sub>3</sub> O <sub>2</sub>	222.1237	4.88
Fenuron	C <sub>9</sub> H <sub>12</sub> N <sub>2</sub> O	165.1022	4.89
Carbendazim	C <sub>9</sub> H <sub>9</sub> N <sub>3</sub> O <sub>2</sub>	192.0768	6.4
2-Hydroxyatrazine	C <sub>8</sub> H <sub>15</sub> N <sub>5</sub> O	198.1349	6.99
Atrazine-Desisopropyl	C <sub>5</sub> H <sub>8</sub> ClN <sub>5</sub>	174.0541	7.21
Metamitron	C <sub>10</sub> H <sub>10</sub> N <sub>4</sub> O	203.0927	7.71
Acetamiprid	C <sub>10</sub> H <sub>11</sub> ClN <sub>4</sub>	223.0745	7.86
Chloridazone	C <sub>10</sub> H <sub>8</sub> ClN <sub>3</sub> O	222.0429	7.88
Crimidine	C <sub>7</sub> H <sub>10</sub> ClN <sub>3</sub>	172.0636	7.89
Pirimicarb	C <sub>11</sub> H <sub>18</sub> N <sub>4</sub> O <sub>2</sub>	239.1503	8.13
Atrazine-Desethyl	C <sub>6</sub> H <sub>10</sub> ClN <sub>5</sub>	188.0697	8.16
Atraton	C <sub>9</sub> H <sub>17</sub> N <sub>5</sub> O	212.1506	8.29
Metoxuron	C <sub>10</sub> H <sub>13</sub> ClN <sub>2</sub> O <sub>2</sub>	229.0738	8.54
2-4-Dimethylphenylformamide	C <sub>9</sub> H <sub>11</sub> NO	150.0913	8.75
Metolcarb	C <sub>9</sub> H <sub>11</sub> NO <sub>2</sub>	166.0863	8.75
Nicosulfuron	C <sub>15</sub> H <sub>18</sub> N <sub>6</sub> O <sub>6</sub> S	411.1081	8.96
Carbofuran	C <sub>12</sub> H <sub>15</sub> NO <sub>3</sub>	222.1125	9.08
Carboxin	C <sub>12</sub> H <sub>13</sub> NO <sub>2</sub> S	236.074	9.26
Fenpropidin	C <sub>19</sub> H <sub>31</sub> N	274.2529	9.32
Fosthiazate	C <sub>9</sub> H <sub>18</sub> NO <sub>3</sub> PS <sub>2</sub>	284.0538	9.45
Cyprazin	C <sub>9</sub> H <sub>14</sub> ClN <sub>5</sub>	228.101	9.69
DEET (diethyltoluamide)	C <sub>12</sub> H <sub>17</sub> NO	192.1383	9.71
Diuron	C <sub>9</sub> H <sub>10</sub> Cl <sub>2</sub> N <sub>2</sub> O	233.0243	9.76
Cycluron	C <sub>11</sub> H <sub>22</sub> N <sub>2</sub> O	199.1805	9.82
Phenmedipham	C <sub>16</sub> H <sub>16</sub> N <sub>2</sub> O <sub>4</sub>	301.1183	9.87
Azoxystrobin	C <sub>22</sub> H <sub>17</sub> N <sub>3</sub> O <sub>5</sub>	404.1241	10.01
Isoxaben	C <sub>18</sub> H <sub>24</sub> N <sub>2</sub> O <sub>4</sub>	333.1809	10.23
Methoxyfenozide	C <sub>22</sub> H <sub>28</sub> N <sub>2</sub> O <sub>3</sub>	369.2173	10.27
Chromafenozide	C <sub>24</sub> H <sub>30</sub> N <sub>2</sub> O <sub>3</sub>	395.2329	10.42
Metolachlor	C <sub>15</sub> H <sub>22</sub> ClNO <sub>2</sub>	284.1412	10.65
Fenothiocarb	C <sub>13</sub> H <sub>19</sub> NO <sub>2</sub> S	254.1209	10.78
Pencycuron	C <sub>19</sub> H <sub>21</sub> ClN <sub>2</sub> O	329.1415	11.05



**Figure S1. Sample <sup>1</sup>H NMR spectrum of urine (a) and LC-MS base peak chromatogram from the same urine sample (b).**



**Figure S2.** Scores plot of the PCA performed on the entire dataset including QC samples. (A) marked by QCs (●) and individual urine samples (□); QC samples form a tight cluster, indicating the analytical reproducibility of the method. (B) Marked by column; no separation by column is visible.

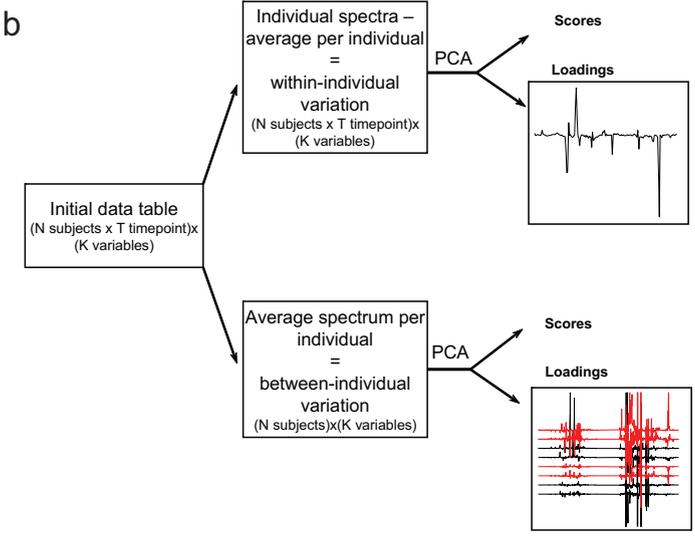
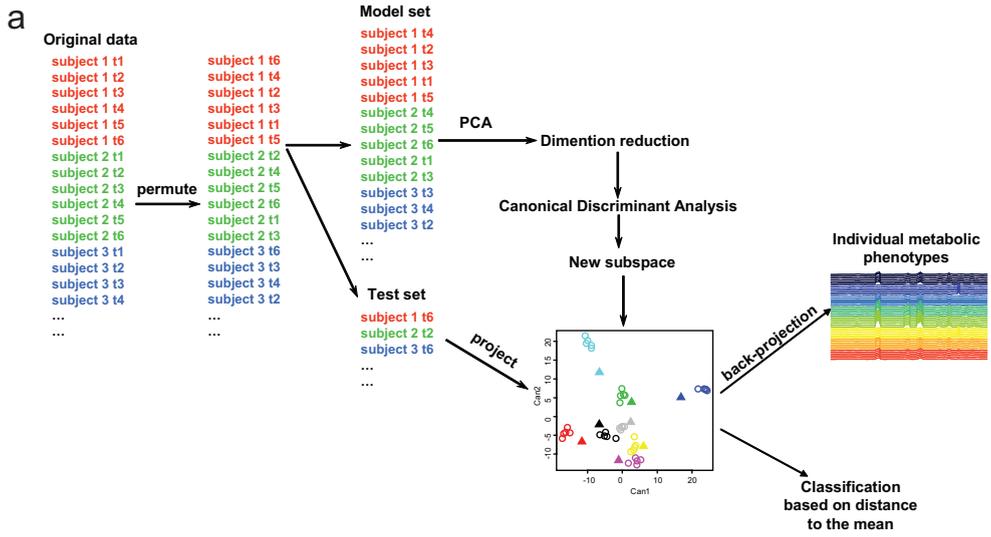
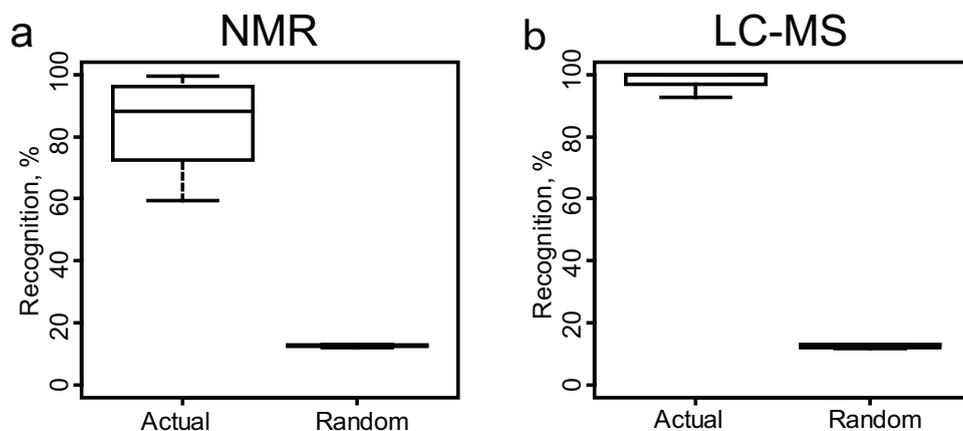
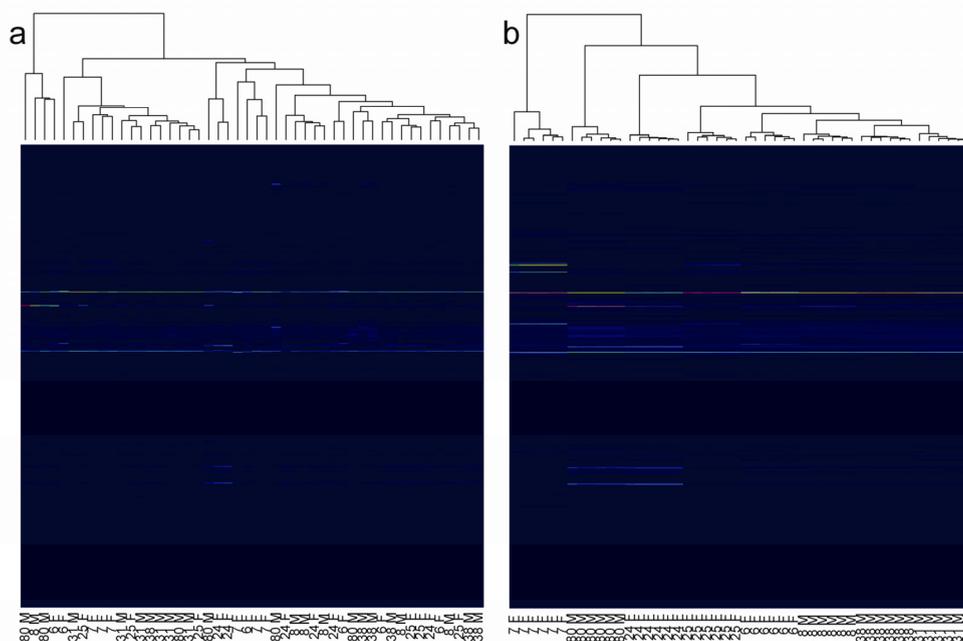


Figure S3. Schematic representation of the two statistical methods used: (a) person recognition, (b) multilevel component analysis.



**Figure S4.** Person recognition based on  $^1\text{H}$  NMR and LC-MS data of urine samples from 8 individuals. (a) Boxplot of the recognition accuracy based on  $^1\text{H}$  NMR spectra for actual (left) and permuted (right) person labels. (b) Boxplot of the recognition accuracy based on LC-MS data for actual (left) and permuted (right) person labels.



**Figure S5.** Heatmap and hierarchical clustering of the initial binned  $^1\text{H}$  NMR table (a) and after back-projection (b). Samples are labeled by their ID and gender.

