

RESEARCH

Open Access



# Bayesian inference for biomarker discovery in proteomics: an analytic solution

Noura Dridi<sup>1,2†</sup>, Audrey Giremus<sup>1†</sup>, Jean-Francois Giovannelli<sup>1\*†</sup> , Caroline Truntzer<sup>3</sup>, Melita Hadzagic<sup>1,4</sup>, Jean-Philippe Charrier<sup>5</sup>, Laurent Gerfault<sup>6,7</sup>, Patrick Ducoroy<sup>3</sup>, Bruno Lacroix<sup>5</sup>, Pierre Grangeat<sup>6,7</sup> and Pascal Roy<sup>8,9,10,11</sup>

## Abstract

This paper addresses the question of biomarker discovery in proteomics. Given clinical data regarding a list of proteins for a set of individuals, the tackled problem is to extract a short subset of proteins the concentrations of which are an indicator of the biological status (healthy or pathological). In this paper, it is formulated as a specific instance of variable selection. The originality is that the proteins are not investigated one after the other but the best partition between discriminant and non-discriminant proteins is directly sought. In this way, correlations between the proteins are intrinsically taken into account in the decision. The developed strategy is derived in a Bayesian setting, and the decision is optimal in the sense that it minimizes a global mean error. It is finally based on the posterior probabilities of the partitions. The main difficulty is to calculate these probabilities since they are based on the so-called evidence that require marginalization of all the unknown model parameters. Two models are presented that relate the status to the protein concentrations, depending whether the latter are biomarkers or not. The first model accounts for biological variabilities by assuming that the concentrations are Gaussian distributed with a mean and a covariance matrix that depend on the status only for the biomarkers. The second one is an extension that also takes into account the technical variabilities that may significantly impact the observed concentrations. The main contributions of the paper are: (1) a new Bayesian formulation of the biomarker selection problem, (2) the closed-form expression of the posterior probabilities in the noiseless case, and (3) a suitable approximated solution in the noisy case. The methods are numerically assessed and compared to the state-of-the-art methods (*t* test, LASSO, Battacharyya distance, FOHSIC) on synthetic and real data from proteins quantified in human serum by mass spectrometry in selected reaction monitoring mode.

**Keywords:** Variable selection, Model selection, Optimal decision, Bayesian approach, Evidence, Hierarchical model, Proteomics, Biomarker

## 1 Introduction

It is now generally recognized that protein expression analysis is crucial in explaining the changes that occur as a part of disease pathogenesis [1, 2]. In this context, recent advances in mass spectrometry (MS) technologies have facilitated the investigation of proteins over a wide range of molecular weights in small biological specimens from blood or urine samples, for instance. Notably, MS

in selected reaction monitoring (SRM) mode has demonstrated its ability to quantify clinical biomarkers in patient sera [3, 4]. Consequently, a large amount of research has been generated in proteomics based on data such as protein mass spectral intensities or protein concentrations obtained from the spectra. Specifically, the focus is on the selection (or discovery) of the “signature profiles,” the so-called biomarkers. They represent, for instance, indicators of normal versus pathogenic biological processes, or positive versus negative pharmacological responses to therapeutic intervention.

Critical to the identification of biomarkers are: (1) the biological variability, i.e., the random variations of the concentrations of proteins between individuals sharing

\*Correspondence: Giova@IMS-Bordeaux.fr

†Equal contributors

<sup>1</sup>IMS (Univ. Bordeaux, CNRS, BINP), 33400 Talence, France

Full list of author information is available at the end of the article

the same biological status [5], and (2) the technical variability, which originates from the imperfections of the measurement process used to obtain the concentrations. Failing to address both of these variabilities within a technique for biomarker identification may significantly impair its performance by resulting in erroneous decision.

Furthermore, since the complexity of a status is unlikely to be manifested through the changes in the characteristics of just one protein, it has generally been acknowledged that a set of proteins should be considered [5–8]. An additional difficulty is that they are possibly correlated, imposing the use of multivariate models to account for all the data simultaneously. These aforementioned issues pose significant challenges in developing efficient and robust statistical techniques for the identification of biomarkers.

The paper tackles the problem of biomarker identification by adopting a Bayesian approach to propose the selection of the optimal set of variables. By providing an elegant and mathematically rigorous framework for incorporating the data and the prior information within a joint probabilistic model, the Bayesian setting allows straightforward modeling of both the technical and the biological variabilities of the data.

The remainder of the paper is organized as follows. Section 2 summarizes the state-of-the-art variable selection methods, discusses their main challenges, and outlines our principal contributions. Section 3 presents the proposed formulation within the Bayesian framework, the proposed models for the data, and the decision strategy. Section 4 describes the data used in the numerical evaluations, together with the results and their analysis. Finally, conclusions are drawn in Section 5. A detailed description of the model and the derivation of the analytic solution is provided in Appendix.

## 2 Related work

The identification of biomarkers for diagnosis or prognosis can be classically formulated as a variable selection problem, and this problem has been paid a lot of attention as a specific instance of model choice. Various methodologies exist that can be broadly classified in two categories: the frequentist hypothesis testing and the Bayesian decision-making.

Frequentist hypothesis testing consists in deciding between two statements, classically referred to as the null and the alternative hypotheses, by comparing a function of the observed data to a threshold. The reader is invited to consult [9] for a comprehensive overview. Two closely related methods have been proposed. On the one hand, Neyman-Pearson tests are designed to ensure the so-called type I error. On the other hand,  $p$  values focus on how strongly the data reject the null hypothesis  $H_0$  by evaluating the probability of obtaining

a value as extreme as the observed one given  $H_0$  is true. In biomarker discovery, a popular approach consists in testing a mean difference between the case and the negative controlled populations using the classical Students'  $t$  test or its variants [7]. The latter is a statistical hypothesis test which indicates whether the difference between two group means most likely reflects that they are samples of two different populations or, on the contrary, is merely explained by the sampling fluctuation. However, as the number of candidate biomarkers increases, multiple hypothesis testing is required resulting in a higher computational cost which may become prohibitive [10]. A first solution is to perform univariate tests for each protein.

This procedure requires an adapted control of the rate of type I errors in this particular setting where multiple hypothesis tests are conducted simultaneously. Two types of procedures were proposed for this purpose, namely the so-called family wise error rate or the false discovery rate [11, 12]. A common criticism of frequentist approaches is that they fail to take into account prior information about the problem at hand such as interdependencies between the different variables.

For a large number of candidate variables, regression analysis [13] provides an alternative to the above-mentioned methods. The principle is to assume that a given outcome is related to a linear combination of a set of explanatory variables called the predictors. In proteomics, logistic regression models are considered that express the probability to have a disease as a function of the protein abundances [14–16]. Then, variable selection is classically performed using stepwise procedures that consists in successively adding or removing predictors, estimating the regression coefficients, and evaluating the goodness of fit of the subsequent model. Different criteria can be considered such as the  $R$ -squared, the adjusted  $R$ -squared, or the Akaike Information Criterion [17, 18]. Such techniques are referred to as backward elimination and forward selection, respectively [13]. However, these selection procedures are prone to overfitting and the variance of the parameter estimates becomes high in the presence of correlated predictors. Regularization methods alleviate these difficulties by considering the minimum of a penalized least squares error as estimate. Since the Ridge regression in 1970 [19], several algorithms have been proposed that differ between one another with respect to the considered penalization of the regression parameters. The well-known LASSO [20] considers a  $L_1$ -norm penalty and has the advantage of directly removing irrelevant predictors by shrinking their coefficients to zero. More recently, the elastic net [21] which combines the advantages of the Ridge and LASSO regressions has been proposed. In the presence of correlated variables, it outperforms LASSO by favoring the selection of sets of variables. A comparison between these methods in application to

genome selection is presented in [22]. Although widely used, regression analysis is based on an ad hoc model that may not reflect the physical nature of the observed data. Further, it does not explicitly accommodate correlations between the candidate biomarkers as well as measurement errors.

The Bayesian framework offers an alternative formulation of model selection. The candidate models are assigned prior probabilities that are combined with the likelihood function to yield the so-called posterior probability. The latter summarizes all the available information to make the decision. In this context, deciding in favor of the a posteriori most probable model is optimal in the sense that it minimizes the risk associated to the 0/1 cost-function. There have been a lot of debates over the use of Bayesian techniques in place of frequentist approaches, but they do not address exactly the same question. Frequentist methods are designed to test the departure of the data from a pre-defined null hypothesis. In contrast, Bayesian selection procedures evaluate the plausibility of a given hypothesis given a set of candidate hypotheses hence are conveniently well-suited to multiple hypotheses testing. Thus, non-nested models can be compared in a straightforward manner. Another fundamental difference is the treatment of unknown model parameters. In the frequentist approach, they are classically replaced by estimates whereas in the Bayesian formulation, they are integrated. The latter procedure has the advantage of automatically penalizing complex models, as discussed in [23], but often leads to intractable calculations. An additional advantage is that correlations between the model variables can be easily accounted for in the design of the prior distributions. As for the integration over the unknown parameters, several solutions have been developed. The Laplace approximation of the integrand leads to the well-known Bayesian information criterion (BIC). As an alternative, numerical integrations can be performed based on stochastic sampling techniques such as Markov Chain Monte Carlo (MCMC) methods [24]. Either across or within model-based techniques can be considered. In the first case, the model index is sampled jointly with the parameters conditionally upon the observations. A well-known algorithm is the Reversible-Jump MCMC but moves between the different parameter spaces are difficult to design. In the second case, posterior samples of the parameters are generated conditionally upon each candidate model and then used to evaluate the integrated likelihood, also called evidence [25]. Nevertheless, the harmonic mean-based estimator exhibits instabilities [26]. Applications of the MCMC Bayesian model selection methods in genomics can be found in [27, 28].

In this paper, a Bayesian setting is adopted to identify a set of protein biomarkers from experimental data consisting of measured protein concentrations and the

associated biological statuses of a population of individuals. The novelty is that the decision is not made protein by protein. As an alternative, the problem is formulated as directly finding the best partition of the list of proteins into two subsets, namely discriminant and non-discriminant, in the sense that it yields the highest posterior probability. Regardless of their discriminative power, the proteins are assumed Gaussian distributed. However, for the subset of biomarkers, the parameters of the Gaussian distribution take different values depending on the biological status whereas this is not the case for the second subset of proteins. The preliminary version of this hierarchical model has been presented in [29]. Its advantages are threefold. First, it is not based on an ad hoc explanatory model unlike regression analysis. Second, the proteins within a given group are assumed a priori correlated and the dependency structure is integrated out along with the remainder of the unknown model parameters so that only the protein partition is estimated. Thus, our approach intrinsically takes into account correlations between the candidate biomarkers. Third, by choosing appropriate conjugate prior distributions for the parameters, the model evidences can be calculated in closed form and there is no need to resort to computationally extensive numerical techniques. Finally, we show that our hierarchical model can be easily extended to address errors in the measured concentrations.

### 3 Problem formulation, proposed models, and methods

To formulate the biomarker selection problem and construct its solution in the proposed framework, we first introduce the basic modeling for the relevant quantities/variables at hand: the biological status, protein concentrations, number of individuals, ... including the descriptions of the considered observation models.

**Distribution for status and concentration** Regarding the biological status, it is denoted by  $b$  and takes two values,  $\mathcal{H}$  and  $\mathcal{P}$ , for healthy and pathological. It is conveniently described by a Bernoulli random variable  $B$  with parameter  $p$

$$B|p \sim \mathcal{B}(b; p). \quad (1)$$

Regarding the proteins, let us note  $P$  is their number and  $\mathbf{x} \in \mathbb{R}^P$  is the collection of their concentrations. Each protein can be discriminant or non-discriminant and then accordingly labeled by  $+$  or  $-$ . The vector  $\mathbf{x}^+$  and  $\mathbf{x}^-$ , with sizes  $P^+$  and  $P^-$  (we have  $P = P^+ + P^-$ ), stand for the respective concentrations. Otherwise, within the  $P$  proteins, there are  $2^P$  possible partitions referred to as  $\delta \in \{+, -\}^P$  since each protein can be discriminant and non-discriminant. Following the clinically observed

behavior, from a probabilistic standpoint, the concentrations are described by normal distributions in order to account for biological variabilities within the populations. Specifically, for the non-discriminant proteins, the concentration vector  $\mathbf{x}^-$  is modeled by a unique multivariate normal distribution with common parameters  $(\mathbf{m}_C, \Gamma_C)$  regardless the biological status:

$$X^- | \mathbf{m}_C, \Gamma_C \sim \mathcal{N}(\mathbf{x}^-; \mathbf{m}_C, \Gamma_C). \quad (2)$$

On the other hand, for discriminant ones, the concentration vector  $\mathbf{x}^+$  is modeled, conditionally on status  $b$ , by a multivariate normal distribution with mean and precision  $(\mathbf{m}_\mathcal{H}, \Gamma_\mathcal{H})$  and  $(\mathbf{m}_\mathcal{P}, \Gamma_\mathcal{P})$  for healthy and pathological, respectively:

$$\begin{cases} X^+ | b = \mathcal{P}, \delta, \mathbf{m}_\mathcal{P}, \Gamma_\mathcal{P} \sim \mathcal{N}(\mathbf{x}^+; \mathbf{m}_\mathcal{P}, \Gamma_\mathcal{P}), \\ X^+ | b = \mathcal{H}, \delta, \mathbf{m}_\mathcal{H}, \Gamma_\mathcal{H} \sim \mathcal{N}(\mathbf{x}^+; \mathbf{m}_\mathcal{H}, \Gamma_\mathcal{H}). \end{cases} \quad (3)$$

Marginally, the concentrations of discriminant proteins are distributed according to a mixture of two Gaussian distributions

$$X^+ | p, \mathbf{m}_\mathcal{P}, \Gamma_\mathcal{P}, \mathbf{m}_\mathcal{H}, \Gamma_\mathcal{H} \sim p \mathcal{N}(\mathbf{x}^+; \mathbf{m}_\mathcal{H}, \Gamma_\mathcal{H}) + (1 - p) \mathcal{N}(\mathbf{x}^+; \mathbf{m}_\mathcal{P}, \Gamma_\mathcal{P}).$$

In addition, it is assumed that  $X^+$  and  $X^-$  are conditionally independent.

The parameters of the distributions are collected in the vector  $\boldsymbol{\theta} = [\mathbf{m}_\mathcal{P}, \Gamma_\mathcal{P}, \mathbf{m}_\mathcal{H}, \Gamma_\mathcal{H}, \mathbf{m}_C, \Gamma_C, p]$  considered as unknown. It is important to keep in mind that the quantity of interest is the partition  $\delta$ .

**Distribution for the individuals** The total number of individual is  $N$  and  $(\mathbf{x}_n, b_n)$  is the  $n$ th concentration vector and status. They are modeled as independent conditionally on  $\boldsymbol{\theta}$ . Let us denote  $\underline{\mathbf{x}}$  (size  $P \times N$ ) as the matrix of concentrations and  $\mathbf{b}$  (size  $N$ ) as the vector of biological statuses. Also, let  $\mathcal{I}_\mathcal{H}$  and  $\mathcal{I}_\mathcal{P}$  be the subsets of indices for healthy and pathological individuals, respectively, and  $N_\mathcal{H}, N_\mathcal{P}$  their cardinality. For notational convenience, we introduce:  $N_C = N_\mathcal{H} + N_\mathcal{P}$  and  $\mathcal{I}_C = \mathcal{I}_\mathcal{P} \cup \mathcal{I}_\mathcal{H}$  (where  $N_C = N$  and  $\mathcal{I}_C = \{1, 2, \dots, N\}$ ).

Given the models (1) for the status and (2)–(3) for the non-discriminant and the discriminant proteins, based on the assumptions that  $X_n^+$  and  $X_n^-$  are conditionally uncorrelated and that the individual concentrations are also conditionally independent, the distribution of the concentrations and status  $(\mathbf{x}, \mathbf{b})$ , given the unknown parameters  $\boldsymbol{\theta}$  and partition  $\delta$ , is:

$$\begin{aligned} f_{\underline{X}, \mathbf{B} | \boldsymbol{\theta}, \Delta}(\mathbf{x}, \mathbf{b} | \boldsymbol{\theta}, \delta) &= \prod_n f_{X, B | \boldsymbol{\theta}, \Delta}(\mathbf{x}_n, b_n | \boldsymbol{\theta}, \delta) \\ &= \prod_n f_{X | B, \boldsymbol{\theta}, \Delta}(\mathbf{x}_n | b_n, \boldsymbol{\theta}, \delta) \mathbb{P}_{B | \boldsymbol{\theta}}(b_n | \boldsymbol{\theta}) \\ &= \prod_{n \in \mathcal{I}_\mathcal{P}} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_\mathcal{P}, \Gamma_\mathcal{P}) \\ &\quad \times \prod_{n \in \mathcal{I}_\mathcal{H}} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_\mathcal{H}, \Gamma_\mathcal{H}) \prod_{n \in \mathcal{I}_C} \mathcal{N}(\mathbf{x}_n^-; \mathbf{m}_C, \Gamma_C) \\ &\quad \times \prod_{n \in \mathcal{I}_C} \mathbb{P}_{B | \boldsymbol{\theta}}(b_n | \boldsymbol{\theta}) \end{aligned} \quad (4)$$

It can be seen (see Appendix 1) that the exponential arguments of the Gaussian distributions in the first three factors can be reformulated based on the empirical means and covariances for each index set  $\mathcal{I}_\mathcal{P}$ ,  $\mathcal{I}_\mathcal{H}$ , and  $\mathcal{I}_C$ . The result is shown here only for the first factor:

$$\begin{aligned} \prod_{n \in \mathcal{I}_\mathcal{P}} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_\mathcal{P}, \Gamma_\mathcal{P}) &= (2\pi)^{-P^+ N_\mathcal{P} / 2} |\Gamma_\mathcal{P}|^{N / 2} \\ &\quad \exp \left[ -\frac{N_\mathcal{P}}{2} \text{Tr} \left( \Gamma_\mathcal{P} \left[ \bar{\mathbf{R}}_\mathcal{P}^+ + (\bar{\mathbf{x}}_\mathcal{P}^+ - \mathbf{m}_\mathcal{P}) (\bar{\mathbf{x}}_\mathcal{P}^+ - \mathbf{m}_\mathcal{P})^\top \right] \right) \right] \end{aligned} \quad (5)$$

where  $\bar{\mathbf{x}}_\mathcal{P}^+$  and  $\bar{\mathbf{R}}_\mathcal{P}^+$  are the empirical mean and covariance of the  $\mathbf{x}_n^+$  for the individuals  $n \in \mathcal{I}_\mathcal{P}$ . Moreover, regarding the probability for the statuses, we have:

$$\prod_{n \in \mathcal{I}_C} \mathbb{P}_{B | \boldsymbol{\theta}}(b_n | \boldsymbol{\theta}) = p^{N_\mathcal{P}} (1 - p)^{N_\mathcal{H}}$$

that is only based on the size of each index set  $\mathcal{I}_\mathcal{P}$  and  $\mathcal{I}_\mathcal{H}$ .

**Observations** Given the previously described concentrations, the proposed developments include two cases for the observation model:

1. In the first one, the concentrations  $\mathbf{x}_n$  are directly observed.
2. The second one accounts for noise: observations write  $\mathbf{y}_n = \mathbf{x}_n + \boldsymbol{\varepsilon}_n$ , where  $\boldsymbol{\varepsilon}_n$  is modeled as a zero-mean Gaussian vector with precision  $\Gamma_\varepsilon$ .

Both of them account for biological variabilities and the latter also includes technological variabilities that arise from both the functioning of the measurement system itself and the post-processing of the spectra. These models are referred to as “noiseless model” and “noisy model”. The corresponding variable selection methods are respectively presented in Sections 3.1 and 3.2.

**Prior distributions** The choice of the prior distribution for the unknown parameters is important. First of all, it must allow us to account for available information

(e.g., nominal values and uncertainties, strong uncertainties,...). Second, it should also enable analytical calculations or numerical computations. To this end, the prior density is chosen as a separable prior distribution, i.e., for the noiseless case

$$\pi_{\Theta|\Delta}(\theta|\delta) = \pi_{\mathcal{P}}(\mathbf{m}_{\mathcal{P}}, \Gamma_{\mathcal{P}}|\delta) \pi_{\mathcal{H}}(\mathbf{m}_{\mathcal{H}}, \Gamma_{\mathcal{H}}|\delta) \pi_{\mathcal{C}}(\mathbf{m}_{\mathcal{C}}, \Gamma_{\mathcal{C}}|\delta) \pi_P(p) \tag{6}$$

for  $(\mathbf{m}_{\mathcal{P}}, \Gamma_{\mathcal{P}})$ ,  $(\mathbf{m}_{\mathcal{H}}, \Gamma_{\mathcal{H}})$ ,  $(\mathbf{m}_{\mathcal{C}}, \Gamma_{\mathcal{C}})$ , and  $p$ . For the noisy case, in addition, we have a factor  $\pi(\Gamma_{\varepsilon})$  for  $\Gamma_{\varepsilon}$  independent from the other variables. Regarding these five variables individually, the choice is driven by a conjugation principle [30]:

- The probability  $p$  is assumed a Beta distributed variable with parameter  $(a, b)$ .
- The  $(\mathbf{m}_{\times}, \Gamma_{\times})$  are assumed to be Normal-Wishart  $\mathcal{NW}$  distributed with parameters  $(\boldsymbol{\mu}_{\times}, \eta_{\times}, \Lambda_{\times}, \nu_{\times})$ , for  $\times \in \{\mathcal{P}, \mathcal{H}, \mathcal{C}\}$ . See Appendix 2.
- The precision  $\Gamma_{\varepsilon}$  is under a Wishart distribution with parameters  $(\Lambda_{\varepsilon}, \nu_{\varepsilon})$ . See Appendix 2.

In the subsequent developments, we proceed with the calculation of the posterior probability for the partitions  $\delta$  in the two cases: noiseless concentrations in Section 3.1 and noisy concentrations in Section 3.2. One of the novelty is an explicit analytical result for the noiseless case and a precise approximation for the noisy case.

### 3.1 Selection using the noiseless data

**Optimal decision-maker** The question of the paper is the one of the identification of a set of discriminant proteins, and it amounts to making a decision regarding the partition  $\delta$ . To build an optimal decision-maker, a binary loss is considered that assigns a null loss to the correct decision and a unitary loss to the incorrect decisions. The risk is the mean loss over the models  $\delta$ , the data  $(\mathbf{x}, \mathbf{b})$ , and the unknown parameters  $\theta$ . The optimal decision-maker is defined as the risk minimizer, and it is known to be the one that selects the most a posteriori probable model. It should be noted that alternative loss functions could be chosen, for instance, one that would penalize differently erroneous partitions depending on the number of biomarkers properly identified. In this case, the decision would still be based on the posterior probabilities but with a different rule. However, our choice not only leads to a simple identification procedure but also naturally prevents overfitting.

Thus, the point is to calculate the posterior probability  $\mathbb{P}_{\Delta|\mathbf{X},\mathbf{B}}(\delta|\mathbf{x}, \mathbf{b})$  for each candidate model  $\delta$ . It is carried out using the Bayes rule as:

$$\mathbb{P}_{\Delta|\mathbf{X},\mathbf{B}}(\delta|\mathbf{x}, \mathbf{b}) = \frac{\mathbb{P}(\Delta = \delta) f_{\mathbf{X},\mathbf{B}|\Delta}(\mathbf{x}, \mathbf{b}|\delta)}{\sum_{\delta} \mathbb{P}(\Delta = \delta) f_{\mathbf{X},\mathbf{B}|\Delta}(\mathbf{x}, \mathbf{b}|\delta)} \tag{7}$$

and it crucially depends on the so-called evidence

$$f_{\mathbf{X},\mathbf{B}|\Delta}(\mathbf{x}, \mathbf{b}|\delta) = \int_{\theta} f_{\mathbf{X},\mathbf{B},\Theta|\Delta}(\mathbf{x}, \mathbf{b}, \theta|\delta) d\theta$$

which can be rewritten by factorization

$$f_{\mathbf{X},\mathbf{B}|\Delta}(\mathbf{x}, \mathbf{b}|\delta) = \int_{\theta} \pi_{\Theta|\Delta}(\theta|\delta) f_{\mathbf{X},\mathbf{B},\Theta,\Delta}(\mathbf{x}, \mathbf{b}|\theta, \delta) d\theta. \tag{8}$$

This calculation is the main difficulty of the paper and more generally in variable and model selection.

In order to carry out this calculation, let us note that the likelihood  $f_{\mathbf{X},\mathbf{B},\Theta,\Delta}(\mathbf{x}, \mathbf{b}|\theta, \delta)$  factorizes (see Eq. (4)) and that the prior  $\pi_{\Theta|\Delta}(\theta|\delta)$  also factorizes (see Eq. (6)). So, we have:

$$f_{\mathbf{X},\mathbf{B}|\Delta}(\mathbf{x}, \mathbf{b}|\delta) = \int_{\theta} \mathcal{NW}(\mathbf{m}_{\mathcal{P}}, \Gamma_{\mathcal{P}}) \mathcal{NW}(\mathbf{m}_{\mathcal{H}}, \Gamma_{\mathcal{H}}) \mathcal{NW}(\mathbf{m}_{\mathcal{C}}, \Gamma_{\mathcal{C}}) \pi_P(p) \prod_{n \in \mathcal{I}_{\mathcal{P}}} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_{\mathcal{P}}, \Gamma_{\mathcal{P}}) \prod_{n \in \mathcal{I}_{\mathcal{H}}} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_{\mathcal{H}}, \Gamma_{\mathcal{H}}) \prod_{n \in \mathcal{I}_{\mathcal{C}}} \mathcal{N}(\mathbf{x}_n^-; \mathbf{m}_{\mathcal{C}}, \Gamma_{\mathcal{C}}) \prod_{n \in \mathcal{I}_{\Theta}} \mathbb{P}_{B|\Theta}(b_n|\theta) d\theta$$

that can itself be factorized in four integrals: three w.r.t. the couple of variable regarding the concentrations  $(\mathbf{m}_{\mathcal{P}}, \Gamma_{\mathcal{P}})$ ,  $(\mathbf{m}_{\mathcal{H}}, \Gamma_{\mathcal{H}})$ ,  $(\mathbf{m}_{\mathcal{C}}, \Gamma_{\mathcal{C}})$  and one w.r.t. the prevalence  $p$ :

$$f_{\mathbf{X},\mathbf{B}|\Delta}(\mathbf{x}, \mathbf{b}|\delta) = \mathcal{I}_{\mathcal{P}}^+(\mathbf{x}) \mathcal{I}_{\mathcal{H}}^+(\mathbf{x}) \mathcal{I}_{\mathcal{C}}^-(\mathbf{x}) \mathcal{J}(\mathbf{b})$$

where the integrals w.r.t. the  $\mathbf{m}_{\times}, \Gamma_{\times}$  read

$$\mathcal{I}_{\times}^{\star}(\mathbf{x}) = \int_{\mathbf{m}_{\times}, \Gamma_{\times}} \mathcal{NW}(\mathbf{m}_{\times}, \Gamma_{\times}) \prod_{n \in \mathcal{I}_{\times}} \mathcal{N}(\mathbf{x}_n^{\star}; \mathbf{m}_{\times}, \Gamma_{\times}) d\mathbf{m}_{\times} d\Gamma_{\times}, \tag{9}$$

with  $\times \in \{\mathcal{P}, \mathcal{H}, \mathcal{C}\}$  and  $\star \in \{+, -\}$  and the integral w.r.t.  $p$  reads

$$\mathcal{J}(\mathbf{b}) = \int_p \pi_P(p) \prod_{n \in \mathcal{I}_{\Theta}} \mathbb{P}_{B|\Theta}(b_n|\theta) dp.$$

As far as the three integrals  $\mathcal{I}_{\times}^{\star}$  are concerned, the calculations are founded on the reduced form (5) for the likelihood (including empirical means and variances) and on the Normal-Wishart prior (23) in Appendix 2 for the  $(\mathbf{m}_{\times}, \Gamma_{\times})$ . Practically, thanks to the conjugacy property,

the integrand in (9) involves the posterior for  $(\mathbf{m}_\times, \Gamma_\times)$  which is Normal-Wishart with parameters

$$\begin{aligned} \nu_\times^{\text{pst}} &= \nu_\times + N_\times \\ \eta_\times^{\text{pst}} &= \eta_\times + N_\times \\ \boldsymbol{\mu}_\times^{\text{pst}} &= (N_\times \bar{\mathbf{x}}_\times + \eta_\times \boldsymbol{\mu}_\times) / (N_\times + \eta_\times) \\ (\boldsymbol{\Lambda}_\times^{\text{pst}})^{-1} &= (\boldsymbol{\Lambda}_\times)^{-1} + N_\times \bar{\mathbf{R}}_\times^* + N_\times \eta_\times (\boldsymbol{\mu}_\times - \bar{\mathbf{x}}_\times^*) \\ &\quad \times (\boldsymbol{\mu}_\times - \bar{\mathbf{x}}_\times^*)^t / (N_\times + \eta_\times) \end{aligned}$$

where ‘‘pst’’ stands for posterior. The finalization of the development relies on the fact that the Normal-Wishart density sums to one: without effective complicate calculus, this yields the result as the ratio of normaliation constants

$$\mathcal{I}_\times^* = \frac{\mathcal{KNW}_\times^{\text{pst}}}{\mathcal{KNW}_\times^{\text{pri}}}$$

where  $\mathcal{KNW}$  is the normalizing constant of the Normal-Wishart density given by (24) in Appendix 2 and where ‘‘pri/pst’’ stands for ‘‘prior/posterior’’.

As a whole, the analytical calculation of the integral in (8) is possible and yields:

$$f_{\underline{\mathbf{X}}, \mathbf{B} | \Delta}(\underline{\mathbf{x}}, \mathbf{b} | \delta) \propto \frac{\mathcal{KNW}_P^{+\text{pst}} \mathcal{KNW}_H^{+\text{pst}} \mathcal{KNW}_C^{-\text{pst}}}{\mathcal{KNW}_P^{+\text{pri}} \mathcal{KNW}_H^{+\text{pri}} \mathcal{KNW}_C^{-\text{pri}}} \quad (10)$$

rendering the usually complex calculations of the evidences straightforward.

Assuming that all candidate models are equally a priori probable, from Eq. (7), the posterior probability across the  $2^P$  models can be inferred. The selected model is the one which maximizes this probability. It should be noted that if prior information is available such as protein-to-protein interactions (PPI’s), it can be taken into account by assigning a higher probability to partitions wherein the related proteins are in the same subset (either discriminant or not). In Eq. (10), the normalizing constants for the posterior distributions depend on the empirical covariance matrices of the population of individuals for the discriminant proteins and the non-discriminant ones, respectively. Their computation is expensive. However, it suffices to compute once the full covariance matrix for all the proteins and then remove the appropriate rows and columns for the  $2^P$  configurations to be tested.

### 3.2 Selection using noisy data

The model presented above assumes that the concentrations are directly observed. Although this assumption leads to closed-form expressions of the posterior probabilities, it may be too simplifying. In practice, the concentrations are known up to an uncertainty and this section extends the above-detailed developments to account for

these uncertainties. However, this comes at the price of intractable calculations, and to overcome this difficulty, we propose a suitable approximation. As introduced above, the measured concentrations are modeled as  $\mathbf{y}_n = \mathbf{x}_n + \boldsymbol{\varepsilon}_n$  where  $\boldsymbol{\varepsilon}_n$  is a zero-mean Gaussian vector with precision  $\Gamma_\varepsilon$ , therefore  $\mathbf{y}_n | \mathbf{x}_n \sim \mathcal{N}(\mathbf{x}_n; \mathbf{0}, \Gamma_\varepsilon)$ . Similarly to the previous section, the vectors of observed concentrations are stacked in a matrix  $\mathbf{y}$  of dimension  $P \times N$ . To select the most probable model, the evidence  $f_{\underline{\mathbf{Y}}, \mathbf{B} | \Delta}(\underline{\mathbf{y}}, \mathbf{b} | \delta)$  must be calculated for each candidate model (it was  $f_{\underline{\mathbf{X}}, \mathbf{B} | \Delta}(\underline{\mathbf{x}}, \mathbf{b} | \delta)$  for the noiseless model). The difficulty is that the calculation of evidence requires not only the marginalization of the model parameters but also of the true concentrations. Furthermore, the precision  $\Gamma_\varepsilon$  is assumed unknown and must also be marginalized. For notational convenience, we state:  $\tilde{\boldsymbol{\theta}} = [\boldsymbol{\theta}, \Gamma_\varepsilon]$  as an extended vector of unknown parameters.

By taking into account the conditional independencies, the evidence can be expressed as:

$$\begin{aligned} f_{\underline{\mathbf{Y}}, \mathbf{B} | \Delta}(\underline{\mathbf{y}}, \mathbf{b} | \delta) &= \int_{\underline{\mathbf{x}}} \int_{\tilde{\boldsymbol{\theta}}} \prod_{n \in \mathcal{I}_C} \mathcal{N}(\mathbf{y}_n; \mathbf{x}_n, \Gamma_\varepsilon) \\ &\quad \prod_{n \in \mathcal{I}_P} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_P, \Gamma_P) \prod_{n \in \mathcal{I}_H} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_H, \Gamma_H) \\ &\quad \prod_{n \in \mathcal{I}_C} \mathcal{N}(\mathbf{x}_n^-; \mathbf{m}_C, \Gamma_C) \prod_{n \in \mathcal{I}_C} \mathbb{P}_{\mathbf{B} | \Theta}(\mathbf{b}_n | \boldsymbol{\theta}) \\ &\quad \pi_{\tilde{\boldsymbol{\theta}}}(\tilde{\boldsymbol{\theta}} | \delta) \, d\tilde{\boldsymbol{\theta}} \, d\underline{\mathbf{x}} \end{aligned} \quad (11)$$

This multiple integral can be handled in several manners. We propose to first perform integration with respect to  $\tilde{\boldsymbol{\theta}}$  and then to integrate the result with respect to  $\underline{\mathbf{x}}$ ; hence, Eq. (11) can be rewritten as:

$$f_{\underline{\mathbf{Y}}, \mathbf{B} | \Delta}(\underline{\mathbf{y}}, \mathbf{b} | \delta) = \int_{\underline{\mathbf{x}}} \mathcal{I}_\varepsilon(\underline{\mathbf{x}}) \mathcal{I}_P^+(\underline{\mathbf{x}}) \mathcal{I}_H^+(\underline{\mathbf{x}}) \mathcal{I}_C^-(\underline{\mathbf{x}}) \, d\underline{\mathbf{x}}$$

where

$$\mathcal{I}_\varepsilon(\underline{\mathbf{x}}) = \int_{\Gamma_\varepsilon} \mathcal{W}(\Gamma_\varepsilon) \prod_{n \in \mathcal{I}_C} \mathcal{N}(\mathbf{y}_n; \mathbf{x}_n, \Gamma_\varepsilon) \, d\Gamma_\varepsilon \quad (12)$$

and in the same manner as previously

$$\begin{aligned} \mathcal{I}_\times^*(\underline{\mathbf{x}}) &= \int_{\mathbf{m}_\times, \Gamma_\times} \mathcal{NW}(\mathbf{m}_\times, \Gamma_\times) \prod_{n \in \mathcal{I}_\times} \mathcal{N}(\mathbf{x}_n^*; \mathbf{m}_\times, \Gamma_\times) \\ &\quad d\mathbf{m}_\times \, d\Gamma_\times \end{aligned} \quad (13)$$

with  $\times \in \{P, H, C\}$  and  $\star \in \{+, -\}$ . The integrals (12) and (13) can be calculated analytically.

On the one hand, the integrand of (12) can be rewritten, up to a proportionality constant, as the distribution of the precision matrix  $\Gamma_\varepsilon$  conditionally upon  $\underline{\mathbf{y}}$  and  $\underline{\mathbf{x}}$ .

This distribution is Wishart with parameters  $(v_\varepsilon^{\text{pst}}, \Lambda_\varepsilon^{\text{pst}})$  expressed as:

$$\begin{cases} v_\varepsilon^{\text{pst}} &= v_\varepsilon^{\text{pri}} + N_C \\ (\Lambda_\varepsilon^{\text{pst}})^{-1} &= (\Lambda_\varepsilon^{\text{pri}})^{-1} + \sum_{n=1}^{N_C} (\mathbf{y}_n - \mathbf{x}_n)(\mathbf{y}_n - \mathbf{x}_n)^t. \end{cases}$$

Then, (12) can be re-arranged as:

$$\begin{aligned} \mathcal{I}_\varepsilon(\underline{\mathbf{x}}) &= \frac{\mathcal{KNW}_\varepsilon^{\text{pst}}}{\mathcal{KNW}_\varepsilon^{\text{pri}}} (2\pi)^{-N_C P/2} \\ &= \mathcal{T}_{P,N} \left( \underline{\mathbf{x}}; v_\varepsilon^{\text{pri}} + 1 - P, \underline{\mathbf{y}}, (\Lambda_\varepsilon^{\text{pri}})^{-1}, \mathbf{I}_N \right) \end{aligned} \quad (14)$$

where  $\mathcal{KNW}_\varepsilon^*$  is the normalization constant of the Wishart distribution and  $\mathcal{T}_{P,N}(\mathbf{T}; q, \mathbf{M}, \Sigma, \Omega)$  denote the matrix t-distribution of parameters  $q, \mathbf{M}, \Sigma$ , and  $\Omega$ , for a matrix  $\mathbf{T}$  of dimensions  $P \times N$ . The expression is recalled in the ‘‘Matrix variate t-distribution’’ section of Appendix 2.

On the other hand, the integrals  $\mathcal{I}_x^*$  can be computed in the same manner as in the previous section using the conjugation property for the couples  $(\mathbf{m}_x, \Gamma_x)$ . Thus, we have:

$$\mathcal{I}_x^*(\underline{\mathbf{x}}) = (2\pi)^{-N_x P^*/2} \frac{\mathcal{KNW}_x^{\text{pst}}}{\mathcal{KNW}_x^{\text{pri}}} \quad (15)$$

with  $\mathcal{KNW}_x^{\text{pst}}$  and  $\mathcal{KNW}_x^{\text{pri}}$  the normalization constants of the prior and posterior Normal-Wishart distributions for  $(\mathbf{m}_x, \Gamma_x)$ , respectively.

In (11), the result of the first integration with respect to  $\tilde{\theta}$  does not yield an expression that can be integrated analytically w.r.t.  $\underline{\mathbf{x}}$ . To address this issue, we propose to take advantage of the fact that a matrix variate t-distribution  $\mathcal{T}_{P,N}(\mathbf{T}; q, \mathbf{M}, \Sigma, \Omega)$  tends to a Gaussian distribution when the degrees of freedom parameter  $q$  tends to infinity.

In a first step, for a high enough value of  $v_\varepsilon^{\text{pri}}$ , (14) can be approximated as:

$$\mathcal{I}_\varepsilon(\underline{\mathbf{x}}) \simeq \prod_{n \in \mathcal{I}_C} \mathcal{N} \left( \mathbf{x}_n; \mathbf{y}_n, \Lambda_\varepsilon^{\text{pri}} \left( v_\varepsilon^{\text{pri}} + 1 - P \right) \right).$$

In a second step, the integrals  $\mathcal{I}_x^*$  can also be approximated by Gaussian distributions although not directly. For this purpose, we focus on the factor of (15) that depends on the true concentration vectors:

$$\begin{aligned} \mathcal{I}_x^*(\underline{\mathbf{x}}) &= C_x^* \left| (\Lambda_x^*)^{-1} + N_x \bar{\mathbf{R}}_x^* + N_x \eta_x (\boldsymbol{\mu}_x - \bar{\mathbf{x}}_x^*) \right. \\ &\quad \left. \times (\boldsymbol{\mu}_x - \bar{\mathbf{x}}_x^*)^t / (N_x + \eta_x) \right|^{-v_x^{\text{pst}}/2} \end{aligned} \quad (16)$$

where  $C_x^*$  is a proportionality constant.

Contrary to (14), the expression (16) does not correspond to a standard probability density function. To make the calculations tractable, we propose to replace the empirical means of the true concentrations  $\bar{\mathbf{x}}_x^*$  by

their approximated values computed from the measured concentrations  $\bar{\mathbf{y}}_x^*$ . By developing the expression of the empirical covariance matrix in (16), it ensues:

$$\begin{aligned} \mathcal{I}_x^*(\underline{\mathbf{x}}) &\simeq C_x^* \left| \Pi_x^* + \sum_{n \in \mathcal{I}_x} (\mathbf{x}_n^* - \bar{\mathbf{y}}_x^*) (\mathbf{x}_n^* - \bar{\mathbf{y}}_x^*)^t \right|^{-v_x^{\text{pst}}/2} \\ &= C_x^* \left| \Pi_x^* \right|^{-v_x^{\text{pst}}/2} \left| \mathbf{I}_P + (\Pi_x^*)^{-1} (\underline{\mathbf{x}}_x^* - \underline{\bar{\mathbf{y}}}_x^*) \right. \\ &\quad \left. \times (\underline{\mathbf{x}}_x^* - \underline{\bar{\mathbf{y}}}_x^*) \right|^{-v_x^{\text{pst}}/2}, \\ &\propto \mathcal{T}_{P,N} \left( \underline{\mathbf{x}}_x^*; v_x^{\text{pst}} + 1 - P^*, \underline{\bar{\mathbf{y}}}_x^*, \Pi_x^*, \mathbf{I}_N \right) \end{aligned}$$

where  $\Pi_x^* = (\Lambda_x^*)^{-1} + N_x \eta_x (\boldsymbol{\mu}_x - \bar{\mathbf{x}}_x^*) (\boldsymbol{\mu}_x - \bar{\mathbf{x}}_x^*)^t (N_x + \eta_x)$  and  $\underline{\bar{\mathbf{y}}}_x^*$  is a matrix of dimensions  $P^* \times N_x$  whose columns are all equal to the empirical mean  $\bar{\mathbf{y}}_x^*$ . Then, provided  $v_x^{\text{pri}}$  is high enough, we can also approximate this matrix variate t-distribution by a Gaussian distribution as for (12):

$$\mathcal{I}_x^*(\underline{\mathbf{x}}) \propto \prod_{n \in \mathcal{I}_x} \mathcal{N} \left( \mathbf{x}_n^*; \bar{\mathbf{y}}_x^*, \Pi_x^* / (v_x^{\text{pri}} - P^* + 1) \right)$$

The performed approximations allow us to express the integrand of the evidence (11) as a product of Gaussian distributions for the true concentration vectors. In this case, the integration can be carried out analytically. By treating separately pathological and healthy individuals, we finally obtain the following expression of the evidence:

$$\begin{aligned} f_{\underline{\mathbf{Y}}|\Delta}(\underline{\mathbf{y}}, \mathbf{b}|\delta) &\simeq \left( \frac{\eta_{\mathcal{P}}^{\text{+pri}} \eta_{\mathcal{H}}^{\text{+pri}}}{\eta_{\mathcal{P}}^{\text{pst}} \eta_{\mathcal{H}}^{\text{pst}}} \right)^{P^+} \left( \frac{\eta_{\mathcal{C}}^{\text{+pri}}}{\eta_{\mathcal{C}}^{\text{pst}}} \right)^{P^-} \\ &\quad \times \left| \Lambda_{\mathcal{P}}^{\text{pri}} \Pi_{\mathcal{P}}^+ \right|^{-v_{\mathcal{P}}^{\text{pri}}} \left| \Lambda_{\mathcal{H}}^{\text{pri}} \Pi_{\mathcal{H}}^+ \right|^{-v_{\mathcal{H}}^{\text{pri}}} \\ &\quad \times \left| \Lambda_{\mathcal{C}}^{\text{pri}} \Pi_{\mathcal{C}}^+ \right|^{-v_{\mathcal{C}}^{\text{pri}}} \\ &\quad \times |\Sigma_{\mathcal{P}}|^{-1/2} |\Sigma_{\mathcal{H}}|^{-1/2} \\ &\quad \exp \left( -\text{Tr} \left( \Sigma_{\mathcal{P}}^{-1} \bar{\mathbf{R}}_{\mathcal{P}}^y + \Sigma_{\mathcal{H}}^{-1} \bar{\mathbf{R}}_{\mathcal{H}}^y \right) / 2 \right) \end{aligned} \quad (17)$$

In this expression,  $\bar{\mathbf{R}}_{\mathcal{P}}^y$  and  $\bar{\mathbf{R}}_{\mathcal{H}}^y$  are the empirical covariance matrices of the measured concentration vectors for the pathological and healthy individuals, respectively. As for  $\Sigma_{\mathcal{H}}$  and  $\Sigma_{\mathcal{P}}$ , they are defined as:

$$\begin{aligned} \Sigma_x &= \frac{\Lambda_\varepsilon^{-1}}{v_\varepsilon^{\text{pri}} + 1 - P} \\ &\quad + \text{blkdiag} \left[ \frac{\Pi_x^+}{v_x^{\text{pri}} + 1 - P^+}, \frac{\Pi_{\mathcal{C}}}{v_{\mathcal{C}}^{\text{pri}} + 1 - P^-} \right] \end{aligned}$$

with  $x \in \{\mathcal{P}, \mathcal{H}\}$  and  $\text{blkdiag}(\mathbf{A}, \mathbf{B})$  the block-diagonal matrix with diagonal elements  $\mathbf{A}$  and  $\mathbf{B}$ .

#### 4 Numerical evaluation

To assess the performance of the proposed method, we have performed extensive numerical experiments using both simulated and real data. They include comparisons with other methods for biomarker identification, namely:

1. The  $t$  test [31], which consists in comparing the means of each protein concentrations between the two cohorts,  $\mathcal{H}$  and  $\mathcal{P}$ . If the null hypothesis, standing for the mean equality, is rejected, then the protein is declared as a biomarker. The type I error, denoted as  $\alpha$ , corresponds to the incorrect rejection of a true null hypothesis. Its value is used to set the  $t$  test decision threshold. In this paper, it is not necessary to adjust the type I errors to account for multivariate effects. The reason is that, for fair comparison purposes, we directly select the setting that leads to the best performance of the test regarding our criterion. This point is commented in Section 4.1 and Fig. 2.
2. The LASSO method [20], based on a linear regression model in which the explanatory variables are the protein concentrations  $\mathbf{x}$ , while the response variables are the biological statuses  $\mathbf{b}$ . The LASSO method estimates the coefficients of the model by minimizing the sum of the squared errors, with a  $L_1$ -norm penalty. Then, a protein is selected as a biomarker if the value of the coefficient corresponding to its concentration is different from zero. This method introduces a regularization parameter denoted  $\lambda$ .
3. The Bhattacharyya distance [32] is a measure of similarity between two probability distributions and by extension between two populations of individuals [32]. For two multivariate normal distributions with respective mean and covariance matrix  $(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$  and  $(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ , it is given by:

$$D_b = \frac{1}{8}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T + \frac{1}{2} \log \left( \frac{\det(\boldsymbol{\Sigma})}{\sqrt{\det(\boldsymbol{\Sigma}_1)\det(\boldsymbol{\Sigma}_2)}} \right)$$

with  $\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)/2$ .

In the sequel, the Bhattacharyya distance is calculated for each protein by replacing the true mean and covariance matrix by their empirical estimates. The protein is declared as discriminant if the distance is greater than a fixed threshold denoted  $t$ . The algorithm is referred to as Bha-distance.

4. The FOHSIC algorithm as introduced in [33]. It performs feature selection based on the Hilbert-Schmidt Independence Criterion (HSIC). The authors propose an unbiased estimator of HSIC and then, assuming the number of significant features is set a priori, use a forward procedure to

select them. In our context, the significant features are the biomarkers.

In this section, we refer to the method from Section 3.1 as the Bayesian Model Selection with Analytical Solution for Noiseless Data (BMS-AS-D) method, while to the method from Section 3.2 as the Bayesian Model Selection with Analytical Solution for Noisy data (BMS-AS-N) method.

Crucial to our approach is the choice of the parameters of the Normal-Wishart densities  $(\nu_\times, \eta_\times, \boldsymbol{\mu}_\times, \boldsymbol{\Lambda}_\times)$  for  $\times \in \{\mathcal{P}, \mathcal{H}, \mathcal{C}\}$ . They are referred to as the hyperparameters since the evidence (10) depends on them. In a non-informative case, the values of these hyperparameters are chosen to be  $(0, 0, 0, \infty)$ , while the proportionality coefficient in (10) has an undetermined form. As an alternative, to tune the parameters, we propose to use a poorly informative prior based on expert knowledge<sup>1</sup> about the corresponding variables (e.g., the values in  $\mu\text{g/ml}$ ). To this end, we take advantage of the expression of the prior mean and the covariance for  $(\mathbf{m}_\times, \boldsymbol{\Gamma}_\times)$  as a function of the hyperparameters:

$$E(\boldsymbol{\Gamma}_\times) = \nu_\times \boldsymbol{\Lambda}_\times \tag{18}$$

$$E(\mathbf{m}_\times) = \boldsymbol{\mu}_\times \tag{19}$$

$$V(\mathbf{m}_\times) = \boldsymbol{\Lambda}_\times^{-1} / [\eta_\times (\nu_\times - P^* - 1)] \tag{20}$$

$$\text{cov}(\boldsymbol{\Gamma}_\times^{ij}, \boldsymbol{\Gamma}_\times^{kl}) = \nu_\times (\boldsymbol{\Lambda}_\times^{il} \boldsymbol{\Lambda}_\times^{jk} + \boldsymbol{\Lambda}_\times^{ik} \boldsymbol{\Lambda}_\times^{jl}), \tag{21}$$

where the superscripts  $i, j$  denote the entry  $(i, j)$  of the matrices,  $E(\cdot)$  and  $V(\cdot)$  refer to the expectation and the covariance matrix of a vector, respectively, while  $\text{cov}(\cdot, \cdot)$  stands for the covariance between two random variables. We also recall that  $*$   $\in \{+, -, "\}$  depending whether the discriminant/non-discriminant subsets of proteins are considered or the whole set. As a consequence, the prior parameters  $(\nu_\times, \eta_\times, \boldsymbol{\mu}_\times, \boldsymbol{\Lambda}_\times)$  can be calculated from (18) to (21) and substituted in (10). Although our choice of prior is not non-informative in the strict sense, it is vague enough so as not to impact biomarker detection. This issue is investigated in the next subsection.

Finally, to calculate (17) in the noisy case, additional hyperparameters for the Wishart probability density function of the noise precision matrix have to be tuned. They are chosen such that  $E(\boldsymbol{\Gamma}_\varepsilon) = \nu_\varepsilon^{\text{pri}} \boldsymbol{\Lambda}_\varepsilon^{\text{pri}}$  and that the elements of the covariance matrix satisfy  $\text{cov}(\boldsymbol{\Gamma}_\varepsilon^{ij}, \boldsymbol{\Gamma}_\varepsilon^{kl}) = \nu_\varepsilon^{\text{pri}} (\boldsymbol{\Lambda}_\varepsilon^{\text{pri},il} \boldsymbol{\Lambda}_\varepsilon^{\text{pri},jk} + \boldsymbol{\Lambda}_\varepsilon^{\text{pri},ik} \boldsymbol{\Lambda}_\varepsilon^{\text{pri},jl})$ . Therefore, by accounting for real-life orders of magnitudes of  $\boldsymbol{\Gamma}_\varepsilon$ , the prior parameters  $(\nu_\varepsilon^{\text{pri}}, \boldsymbol{\Lambda}_\varepsilon^{\text{pri}})$  can be calculated and substituted in the probability (10).

In the next sections, we present the results of the numerical evaluations of the proposed methods using both simulated and real data.



#### 4.1 Evaluation using simulated data

##### 4.1.1 Description of the simulated data and performance index

We consider the concentrations of a list of  $P$  proteins for a group which comprises  $N_{\mathcal{H}}$  healthy and  $N_{\mathcal{P}}$  pathological individuals, respectively, with  $N_{\mathcal{H}} + N_{\mathcal{P}} = N$ . The possible partitions for discriminant/non-discriminant proteins thus amount to  $2^P$  and they are referred to as true models. For each true model,  $N_r = 10^5$  data realizations are simulated, hence the total number of realizations equals  $N_r 2^P$ .

On the one hand, the noiseless data comprise the biological statuses  $b_n$  and the actual protein concentrations  $\mathbf{x}_n$  of the  $N$  individuals and are generated as follows. The biological statuses are sampled from the Bernoulli distribution of parameter  $p$ , where  $p$  is assumed Beta distributed of parameters  $a = 1$  and  $b = 1$ , which corresponds to a uniform distribution. The concentrations of the subset of discriminant proteins are generated from the Gaussian distributions,  $\mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_{\mathcal{H}}, \Gamma_{\mathcal{H}})$  or  $\mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_{\mathcal{P}}, \Gamma_{\mathcal{P}})$ , depending on the simulated biological status. The subset of non-discriminant proteins are sampled from  $\mathcal{N}(\mathbf{x}_n^-; \mathbf{m}_{\mathcal{C}}, \Gamma_{\mathcal{C}})$ . The parameters  $(\mathbf{m}_{\times}, \Gamma_{\times})$ , where  $\times \in \{\mathcal{P}, \mathcal{H}, \mathcal{C}\}$ , are distributed according to the Normal-Wishart distribution  $\mathcal{NW}(v_{\times}, \eta_{\times}, \boldsymbol{\mu}_{\times}, \boldsymbol{\Lambda}_{\times})$ . The orders of magnitudes for  $(\mathbf{m}_{\times}, \Gamma_{\times})$  are specified as:  $E(\mathbf{m}_{\times}) = 10^3 \mathbf{1}_{P^*}$ ,  $V(\mathbf{m}_{\times}) = 10^4 \mathbf{1}_{P^*}$ ,  $E(\Gamma_{\times}) = 10^3 \mathbf{I}_{P^*}$ ,  $V(\Gamma_{\times}) = 10^4 \mathbf{I}_{P^*}$ , where  $\mathbf{1}_{P^*}$  denotes a vector of size  $P^*$  whose elements are all equal to 1 and  $\mathbf{I}_{P^*}$  is the identity matrix of size  $P^*$ . The same order of magnitude is considered for healthy, pathological, and common parameters, that is to say  $\times \in \{\mathcal{H}, \mathcal{P}, \mathcal{C}\}$ . These a priori information are used to tune the hyperparameters as given by (18)–(21).

On the other hand, the noisy data include the biological statuses  $b_n$  and the observed protein concentrations  $\mathbf{y}_n$  for  $n = 1, 2, \dots, N$ . The protein concentrations  $\mathbf{x}_n$  are generated as in the case of the noiseless observations by using the same hyperparameter setting. As for the noise  $\boldsymbol{\varepsilon}_n$ , it is sampled as a zero-mean multivariate Gaussian random vector with precision matrix  $\Gamma_{\varepsilon}$ . The latter is generated from a Wishart density with parameters  $(v_{\varepsilon}^{\text{pri}}, \boldsymbol{\Lambda}_{\varepsilon}^{\text{pri}})$ . In order to determine these hyperparameters, the a priori information is specified as:  $E(\Gamma_{\varepsilon}) = 10^{-2} \mathbf{I}_P$ ,  $V(\Gamma_{\varepsilon}) = 10^{-5} \mathbf{I}_P$ . Note here that  $\Gamma_{\varepsilon}$  measures the precision (inverse variance), thus the lower  $\Gamma_{\varepsilon}$  is, the stronger the noise is.

For each data set, the posterior probability is computed for all possible partitions according to (10) or (17) for the BMS-AS-D and BMS-AS-N methods, respectively. Then, the most probable partition is selected.

The performance is measured in terms of the error rate  $\tau$ , defined as:

$$\tau(\%) = \left( \sum_{i=1}^{2^P} E_i / (N_r \times 2^P) \right) \times 100$$

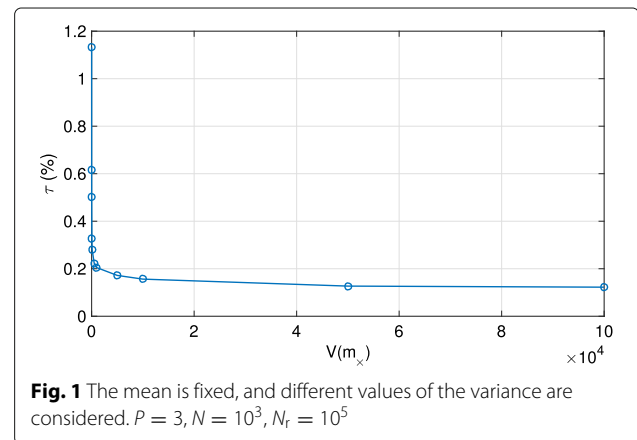
where  $E_i$  is the number of realizations of the  $i$ th partition for which the selected model is different from the true one. It should be noted that the usual type I and type II errors apply when the biomarker identification is made protein after protein but they are not relevant here, since the proteins are addressed as a whole. Indeed, as underlined in Section 3, the proposed Bayesian formulation relies on a different paradigm than the existing methods: any wrong or correct decision regards the list of proteins as a whole (and not protein-wise). Furthermore, the proposed approach minimizes a global risk and cannot be directly related to a given false discovery rate. It ensues that  $\tau$ , which encompasses all types of errors, appears as a suitable criterion to measure the performance.

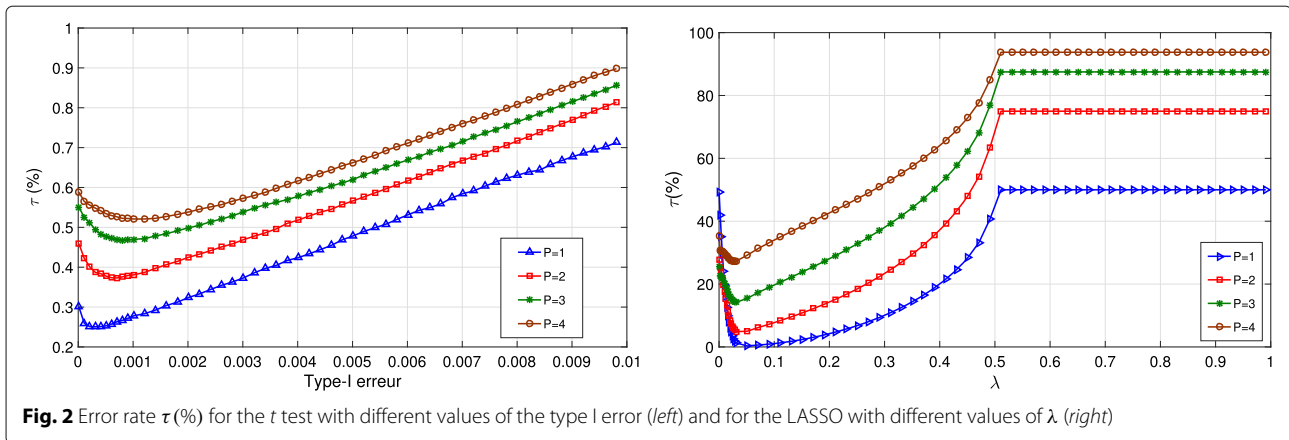
##### 4.1.2 Results for the noiseless model

First of all, we investigate the sensitivity of the error rate  $\tau(\%)$  to the hyperparameter tuning. In Fig. 1, we plot the latter as a function of the considered variance  $V(\mathbf{m}_{\times})$  in Eqs. (18)–(21). We can observe that provided the variance is chosen superior to a given threshold, it has no impact on the biomarker detection performance. This is the case for our settings.

Before going further in analyzing the performance, it should be noted that the  $t$  test, the LASSO, and the Bat-tacharyya distance all require the setting of a parameter: the type I error  $\alpha$ , the regularization parameter  $\lambda$ , and the threshold  $t$ , respectively. So as not to favor our approach, we have run all the algorithms for different values of these parameters and we have selected the best one (in order to get the lowest error rate). Such a procedure cannot be applied on real data, but it allows us to compare the proposed method to the best version of the alternative approaches. The results are given in Fig. 2 for the  $t$  test and the LASSO.

The performance of the BMS-AS-D method has first been evaluated as a function of the number of proteins and results are shown in Table 1. The superior performance





of the BMS-AS-D method with respect to the  $t$  test is expected since the BMS-AS-D makes a decision jointly across all the proteins while accounting for all possible correlations between them. This issue is not addressed within the  $t$  test which is univariate, i.e., each protein is tested separately. The same observation can be made for the Bhattacharyya distance. Indeed, the error rate of the  $t$  test and the Bhattacharyya distance are very close. As for the superiority of the BMS-AS-D method over the LASSO, it is due to the fact that the latter makes the assumption of an arbitrary linear regression relationship between the biological statuses and the protein concentrations. Moreover, the difference in performance increases with the number of proteins, since the possibility of correlation increases. Indeed, in the presence of correlated variables, the number of significant variables is known to be over-estimated with the LASSO algorithm. This fact confirms the relevance of the BMS-AS-D method which can accommodate correlations between the variables.

Regarding the number of individuals, Table 2 shows that, as expected, the performance of all the methods improves as the number of considered individuals increases. In particular, the better performance of the BMS-AS-D method is explained by the reduced variance of the protein concentration posterior distribution for a large number of individuals. Furthermore, the BMS-AS-D method outperforms the  $t$  test, the LASSO and the Bhattacharyya distance, regardless the number of individuals. Finally, even if the number of configurations to be

tested increases exponentially with the number of proteins, the computational cost is kept reasonable owing to the analytical expression of the posterior probabilities for the different partitions.

Last but not the least, our method can also be run for a fixed number of biomarkers as it is the case of many current feature identification algorithms such as the FOHSIC. Only the partitions with the proper number of biomarkers are studied, which amounts to assigning a null prior probability to the others. When the number of biomarkers is chosen as  $M \leq P$ , the number of posterior probabilities to compute is limited to  $C_P^M$  instead of  $2^P$ . Under this assumption, the performance of the BMS-AS-D is compared to that of the FOHSIC in the Tables 3, 4, 5 and 6. To run the algorithm with large  $P (\geq 8)$ , the number of realizations is reduced to  $10^3$ .

As shown in Tables 3, 5, and 6, the BMS-AS-D algorithm outperforms the FOHSIC one. This is explained by the fact that the BMS-AS-D algorithm makes a multivariate decision on the whole set of proteins, while the FOHSIC uses a forward procedure which can lead to error accumulation. Indeed, if any detected biomarker in the sequence is false, then the final selected model is bound to be erroneous. Furthermore, the BMS-AS-D is also faster than the FOHSIC, as illustrated in Table 4. More precisely, Table 5 shows the error rate  $\tau$  for the FOHSIC and the Bayesian algorithm for  $P = 8$  and different number of biomarkers. Conversely, the results proposed in Table 6 are obtained with the number of biomarkers fixed to  $M = 4$  while

**Table 1** Noiseless data:  $\tau$  (%) for different value of  $P, N = 1000$

$P$	1	2	3	4
Best $t$ test	0.2555	0.4025	0.471	0.5209
Best LASSO	0.3185	4.865	14.207	27.2655
Best Bha-distance	0.2245	0.3743	0.4496	0.5008
BMS-AS-D	0.0935	0.1434	0.1563	0.1616

**Table 2** Noiseless data:  $\tau$  (%) for different values of  $N, P = 3$

$N$	100	500	1000
$t$ test	4.041	0.8938	0.471
LASSO	20.5541	15.9870	14.207
Best Bha-distance	3.6970	0.8194	0.4496
BMS-AS-D	1.71	0.32	0.1563

**Table 3**  $\tau$  (%)  $N = 1000$  and  $P = 3$ 

Number of biomarkers	2
BMS-AS-D	0.0723
FOHSIC	0.2907

the number of proteins  $P$  is varied. As expected, the performance of the FOHSIC algorithm is degraded when increasing the number of proteins while the opposite is observed for the BMS-AS-D. Thus, even for large  $P$ , the Bayesian algorithm outperforms the FOHSIC.

#### 4.1.3 Results for the noisy model

The performance of the BMS-AS-N and the BMS-AS-D algorithms is first studied as a function of the number of proteins  $P$  and the number of individuals  $N$ , for a fixed noise level. Then, the BMS-AS-N and the BMS-AS-D methods are compared for different noise conditions.

Table 7 reports the error rate  $\tau$  (%) for the BMS-AS-N and the BMS-AS-D methods. The value of the error rate for the BMS-AS-D method is increased as compared to the results given in Table 1. This is due to the fact that the BMS-AS-D relies on a noiseless model, i.e., it processes the noisy data as if they were noiseless protein concentrations. That is why this result is expected, especially given the specified severe noise conditions ( $E(\Gamma_\varepsilon) = 10^{-2} \mathbf{I}_P$ ). As a consequence, the performance of the BMS-AS-N method is significantly better than the BMS-AS-D one. Also, the difference in performance increases with the number of proteins, since for large sets of proteins, the number of candidate models increases and the estimation becomes more difficult.

The performance of the BMS-AS-N method is also evaluated for several numbers of individuals:  $N = 100, 500$ , and  $1000$ . Table 8 shows the results where it can be observed that for  $N = 100$ , the error rate is higher; however, it does not exceed 13%. For the BMS-AS-D method, the results are even poorer because of the impact of the noise on smaller sample sizes.

To assess the importance of taking into account the noise in the model, the respective performances of the BMS-AS-D and the BMS-AS-N methods are also compared for different noise levels. We recall that the former is designed from the noiseless data model, while the latter specifically addresses the noisy data model. The noise power is measured by the mean value of the noise variance  $E(\Gamma_\varepsilon)$ , which is varied in the simulations. Table 9

**Table 4** Execution time for one simulation  $N = 1000$  and  $P = 3$ 

Number of biomarkers	2
BMS-AS-D	0.202 s
FOHSIC	0.732 s

**Table 5**  $\tau$  (%) for  $P = 8$ ,  $N_r = 10^3$ , and  $N = 500$ 

$M$ :	4	6
BMS-AS-D	0.3014	0.2107
FOHSIC	0.8271	0.9107

shows the error rate for both methods as a function of this parameter: The BMS-AS-N method always outperforms the BMS-AS-D one. In the absence of noise, the BMS-AS-N method becomes equivalent to the optimal noiseless method. These results confirm the relevance of the method, especially for high noise levels.

As a conclusion, the results confirm the good performance of the proposed BMS-AS-N method which is also not too computationally intensive by means of the analytical approximation of the posterior probabilities.

#### 4.2 Evaluation using the real data

The primary goal of this paper was to present a novel methodology for biomarker identification that relaxes classical simplifying assumptions on the data model and then to evaluate it on simulated data. Nevertheless, we had at our disposal a batch of real data<sup>2</sup> and we used it to carry out a preliminary study of the BMS-AS-N method. The data are composed of 206 samples: 105 with the status  $\mathcal{H}$  (including 76 patients from blood donors and 29 with negative colonoscopy), 101 with malignant tumor<sup>3</sup>, i.e. with status  $\mathcal{P}$ . The latter are structured as follows: 24 patients in the 'stage one' of the cancer, 26 patients in the 'stage two', 23 patients in the 'stage three', 25 patients in the 'stage four', and three missing values. The protein concentrations are obtained using the Bayesian inversion method developed in [34] from measurements of SRM spectra according to the methodology described in [35]. For each sample, the concentrations of 21 proteins are measured (14-3-3 protein sigma; 78-kDa glucose-regulated protein; protein S100-A11; calmodulin; calreticulin; peptidyl-prolyl cis-trans isomerase A; defensin-5; defensin-6; heat shock cognate 71 kDa protein; fatty acid-binding protein, intestinal; fatty acid-binding protein, liver (LFABP); stress-70 protein, mitochondrial; protein disulfide-isomerase (PDI); protein disulfide-isomerase A6 (PDIA6); phosphoglycerate kinase 1; retinol-binding protein 4; peroxiredoxin 5, mitochondrial; protein S100-A14; triosephosphate isomerase; villin-1 (Villin); Vimentin). Only one of the proteins in the sample, named LFABP, was previously identified by SRM as a biomarker [36]. To

**Table 6**  $\tau$  (%)  $M = 4$ ,  $N_r = 10^3$ , and  $N = 500$ 

$P$	8	12
BMS-AS-D	0.3014	0.1818
FOHSIC	0.8271	0.8103

**Table 7** Noisy data:  $\tau$  (%) for different values of  $P$ ,  $N = 1000$ 

$P$	1	2	3	4
BMS-AS-D	16.82	32.784	48.740	63.686
BMS-AS-N	1.38	2.682	4.158	5.698

calculate the hyperparameters (18)–(21), empirical orders of magnitudes for  $(\mathbf{m}_\times, \Gamma_\times)$  (e.g.,  $\mu\text{g/ml}$ ) are used as specified:  $E(\mathbf{m}_\times) = 10^2 \mathbf{I}_{P^*}$ ,  $V(\mathbf{m}_\times) = 10^3 \mathbf{I}_{P^*}$ ,  $E(\Gamma_\times) = 10^3 \mathbf{I}_{P^*}$ ,  $V(\Gamma_\times) = 10 \mathbf{I}_{P^*}$ .

For this data set, the posterior probability is computed for each of the  $2^{21}$  possible partitions according to (17) for the BMS-AS-N methods. Table 10 presents the four most probable partitions, with their probabilities. By far, the most probable partition (probability 0.9986) is: LFAPB is discriminant and the remaining 20 proteins are non-discriminant. The second most probable partition is: the whole set of protein is non-discriminant (probability 0.001361). The third and the fourth ones select two discriminant proteins and the 19 other are non-discriminant: (LFABP, Villin) and (LFABP, PDIA6), with probability smaller than  $10^{-6}$ . As a conclusion, the proteins are all declared as non-discriminant, except the LFAPB. This study confirms that our method correctly identifies the valid biomarker.

Despite the large number of models to compare (about two millions candidate models), the computation time is just 1 h. This short computation time is made possible by the analytical calculation of the posterior probability, avoiding the use of extensive numerical integration methods such as for instance MCMC algorithms [24].

## 5 Synthesis and perspectives

Biomarker discovery is a challenging task of the utmost interest for the diagnosis and prognosis of diseases. This paper presents a statistical approach based on variable selection. It is developed in a Bayesian framework that relies on an optimal strategy, i.e., the minimization of an error risk. Given  $P$  candidate proteins, the proposed procedure compares the probability of the  $2^P$  partitions (subset of discriminant versus subset of non-discriminant proteins). The most a posteriori probable partition is finally retained and thus defines the selected variables. The main difficulty is the required integration with respect to all the unknown model parameters. An important contribution is to provide a closed-form expression of the probabilities for noiseless observations

**Table 8** Noisy data:  $\tau$  (%) for different valued of  $N$ ,  $P = 3$ 

$N$	100	500	1000
BMS-AS-D	86.850	71.615	48.740
BMS-AS-N	12.477	5.781	4.158

**Table 9**  $\tau$  (%) for  $N = 1000$ ,  $P = 3$ 

$E(\Gamma_\epsilon)(\times \mathbf{I}_P)$	$10^{-2}$	$10^{-1}$	1	10	$10^2$
BMS-AS-D	48.740	12.304	3.064	0.813	0.291
BMS-AS-N	4.158	1.395	0.567	0.303	0.220

and a sensible approximation for noisy observations. The proposed method proved to be well-suited for variable selection in a complex context. Its effectiveness is assessed by a theoretical characterization and numerical studies (on simulated and real data) which are in accordance with the theoretical optimality. Furthermore, the proposed method compares favorably with the  $t$  test, the LASSO, the Battacharrya distance, and the FOHSIC.

From a methodological standpoint, several perspectives can be considered. Regarding the concentrations, non-Gaussian distributions, e.g., Gamma or Wishart models, could be relevant. Regarding the status, a possible development could account for possible errors in the given status. In this case, an additional level should be appended to the hierarchical Bayesian model. It would include a prior probability for an erroneous status.

As for the applicative perspectives, we plan to further take advantage of the performance of the method in other clinical data sets or in other biomedical fields (e.g., genomics, transcriptomics...). In addition, we also intend to make use of the method in other domains, for instance, in astrophysics (identification of pertinent features in order to classify galaxies), or for complex structures and industrial processes (identification of indicators for detection and diagnosis of damages or faults, analysis of fatigue and aging prevention,...).

## Endnotes

<sup>1</sup> The knowledge about orders of magnitudes of the concentration values is acquired from the real data set provided by bioMérieux (Technology Research Department), France.

<sup>2</sup> SRM measurements provided by bioMérieux (Technology Research Department), France

<sup>3</sup> colorectal cancer

## Appendix 1

### Reduction of the concentration distribution

This section explains how the exponential arguments of the Gaussian distributions in (4) can be reformulated

**Table 10** The four most probable partitions, for real data  $P = 21$ 

Declared biomarker	LFABP	No biomarker	LFABP and Villin	LFABP and PDIA6
Probability	9.986 $\times 10^{-1}$	1.361 $\times 10^{-3}$	9.762 $\times 10^{-7}$	8.297 $\times 10^{-7}$

based on the empirical means and covariances of the concentrations, yielding relation (5). We have:

$$\prod_{n \in \mathcal{I}_P} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_P, \Gamma_P) = \prod_{n \in \mathcal{I}_P} (2\pi)^{-P^+/2} |\Gamma_P|^{1/2} \exp\left[-(\mathbf{x}_n - \mathbf{m}_P)^t \Gamma_P (\mathbf{x}_n - \mathbf{m}_P) / 2\right] = (2\pi)^{-P^+ N_P / 2} |\Gamma_P|^{N_P / 2} \exp\left[-\sum_{n \in \mathcal{I}_P} (\mathbf{x}_n - \mathbf{m}_P)^t \Gamma_P (\mathbf{x}_n - \mathbf{m}_P) / 2\right]$$

The idea is to re-arrange the sum in the exponential as a function of the empirical mean  $\bar{\mathbf{x}}_P^+$  and the empirical variance  $\bar{\mathbf{R}}_P^+$ . To this end, for the sake of calculation simplicity, we can write:  $\mathbf{x}_n - \mathbf{m}_P = (\mathbf{x}_n - \bar{\mathbf{x}}_P^+) - (\bar{\mathbf{x}}_P^+ - \mathbf{m}_P)$  and develop the sum of product. Then, using the fact that  $\mathbf{u}^t \mathbf{M} \mathbf{u} = \text{Tr}(\mathbf{u}^t \mathbf{M} \mathbf{u}) = \text{Tr}(\mathbf{M} \mathbf{u} \mathbf{u}^t)$ , the product is written as function of empirical mean and variance

$$\prod_{n \in \mathcal{I}_P} \mathcal{N}(\mathbf{x}_n^+; \mathbf{m}_P, \Gamma_P) = (2\pi)^{-P^+ N_P / 2} |\Gamma_P|^{N_P / 2} \exp\left[-\frac{N_P}{2} \text{Tr}\left(\Gamma_P \left[\bar{\mathbf{R}}_P^+ + (\bar{\mathbf{x}}_P^+ - \mathbf{m}_P)(\bar{\mathbf{x}}_P^+ - \mathbf{m}_P)^t\right]\right)\right] \tag{22}$$

which allows easier handling of the Normal-Wishart prior.

## Appendix 2

### Wishart, Normal-Wishart, and matrix variate t-distribution Wishart

The Wishart density probability function for a  $P \times P$  matrix  $\Gamma$  is driven by two parameters: a degree of freedom  $\nu$  (real and larger than  $P - 1$ ) and a matrix  $\Lambda$  (positive and definite) referred to as the scale matrix. It reads

$$\mathcal{W}(\Gamma; \Lambda, \nu) = \mathcal{KW}^{-1} \det[\Gamma]^{(\nu-P-1)/2} \exp\left[-\text{Tr}[\Gamma \Lambda^{-1}] / 2\right]$$

The normalizing constant  $\mathcal{KW}$  depends on  $\nu$  and  $\Lambda$ :

$$\mathcal{KW} = \mathcal{KW}(\nu, \Lambda) = 2^{\nu P / 2} \det[\Lambda]^{P / 2} \gamma_P(\nu / 2)$$

where  $\gamma_P$  is the  $P$ -dimensional Gamma function.

### Normal-Wishart

For a couple  $(\mathbf{m}, \Gamma)$ , where  $\mathbf{m}$  is a  $P$ -dimensional vector and  $\Gamma$  a  $P \times P$  matrix, the Normal-Wishart density is controlled by four parameters  $\nu, \eta, \boldsymbol{\mu}, \Lambda$ . It reads:

$$\mathcal{NW}(\mathbf{m}, \Gamma; \nu, \eta, \boldsymbol{\mu}, \Lambda) = \mathcal{KNW}^{-1} \det[\Gamma]^{(\nu-P)/2} \exp\left[-\left[\text{Tr}[\Gamma \Lambda^{-1}] + \eta(\mathbf{m} - \boldsymbol{\mu})^t \Gamma (\mathbf{m} - \boldsymbol{\mu})\right] / 2\right] \tag{23}$$

and the normalizing constant is:

$$\mathcal{KNW} = \mathcal{KNW}(\nu, \eta, \boldsymbol{\mu}, \Lambda) = (2\pi)^{P/2} 2^{\nu P / 2} \eta^{-P/2} \det[\Lambda]^{P/2} \gamma_P(\nu/2), \tag{24}$$

and it does not depend on  $\boldsymbol{\mu}$  (that is a position parameter).

### Matrix variate t-distribution

The random matrix  $T$  of dimension  $(P \times N)$  is said to have a matrix variate t-distribution with parameters  $\mathbf{M}, \boldsymbol{\Sigma}, \boldsymbol{\Omega}$ , and  $q$  if its probability density function is given by

$$\mathcal{T}_{P,N}(q, \mathbf{M}, q\boldsymbol{\Sigma}, \boldsymbol{\Omega}) = \frac{\gamma_P\left(\left[\frac{q+N+P-1}{2}\right]\right)}{\gamma_P\left(\left[\frac{P+p-1}{2}\right]\right)} \pi^{-NP/2} |\boldsymbol{\Sigma}|^{-N/2} |\boldsymbol{\Omega}|^{-P/2} |\mathbf{I}_P + \boldsymbol{\Sigma}^{-1}(\mathbf{T} - \mathbf{M})\boldsymbol{\Omega}^{-1}(\mathbf{T} - \mathbf{M})^t|^{-[q+N+P-1]/2}$$

where  $T$  and  $M$  are  $P \times N$  matrices,  $\boldsymbol{\Omega}$  and  $\boldsymbol{\Sigma}$  are positive-definite matrices with respective sizes  $N \times N$  and  $P \times P$  and  $q > 0$ .

When  $q$  tends to infinity, the distribution of  $T$  tends to a Gaussian distribution with mean  $M$  and covariance  $\boldsymbol{\Sigma} \otimes \boldsymbol{\Omega}$  that is to say  $\mathcal{N}(T; M, \boldsymbol{\Sigma} \otimes \boldsymbol{\Omega})$ , where  $\otimes$  is the Kronecker product.

### Funding

This work was supported by the French National Research Agency through the Bayesian Hierarchical Inversion in Proteomics Project under Contract ANR-2010-BLAN-0313.

### Authors' contributions

JFG proposed the initial idea of the method and provided the first theoretical developments. He also importantly contributed to the writing of the paper. AG proposed the extension to noisy data and importantly contributed to the writing of the paper. ND provided the largest part of the Matlab developments and numerical assessment. She also contributed to the writing of the paper. CT provided input in the proteomics fields. She contributed to the comparison with existing methods. She also contributed to the manuscript and proposed valuable comments. MH contributed to the Matlab development and the numerical assessment. She also contributed to the writing of the paper. JPC provided input in the SRM field and with the result interpretation. He acquired and provided real SRM spectra to test the algorithm. He has revised the manuscript. LG has developed the BHI processing algorithms for the MRM mode. He has computed the protein concentration profiles for the MRM clinical data set used for the evaluation on real data. He contributed to the result interpretation. PD provided expertise regarding proteomics. BL contributed to the early developments of the BHI-PRO project. PG was the BHI-PRO project manager. He has coordinated the conception of the processing algorithms and the interpretation of the results. He has revised the manuscript. PR coordinated biostatistical developments and the interpretation of the results. He has revised the manuscript. All authors read and approved the final manuscript.

### Competing interests

JPC and BL are employed by bioMérieux. The other authors declare that they have no competing interests.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Author details**

<sup>1</sup>IMS (Univ. Bordeaux, CNRS, BINP), 33400 Talence, France. <sup>2</sup>National Engineering School of Gabes (ENIG), University of Gabes, Gabes, Tunisia. <sup>3</sup>CLIPP, Pôle de Recherche Université de Bourgogne, 21000 Dijon, France. <sup>4</sup>NATO STO Centre for Maritime Research and Experimentation, 19126 La Spezia, Italy. <sup>5</sup>Technology Research Department, Innovation Unit, bioMérieux SA, Marcy l'Étoile, France. <sup>6</sup>Univ. Grenoble Alpes, F-38000 Grenoble, France. <sup>7</sup>CEA, LETI, MINATEC Campus, F-38054 Grenoble, France. <sup>8</sup>Service de Biostatistique - Bioinformatique, Hospices Civils de Lyon, Lyon, France. <sup>9</sup>CNRS UMR 5558, LBBE, Équipe Biostatistique Santé, Villeurbanne, France. <sup>10</sup>Université de Lyon, Université Claude Bernard Lyon 1, Lyon, France. <sup>11</sup>Pôle Rhône-Alpes de Bioinformatique, Université Claude Bernard - Lyon 1, 69622 Villeurbanne, France.

Received: 4 August 2016 Accepted: 21 June 2017

Published online: 14 July 2017

**References**

1. S Srivastava, *Informatics in Proteomics*, ser. Statistics: a series of textbooks and monographs. (CRC Press, Boca Raton, 2005)
2. KA Do, P Muller, M Vannucci, *Bayesian inference for gene expression and proteomics*. (Cambridge University Press, Cambridge, England, 2006)
3. T Fortin, A Salvador, JP Charrier, C Lenz, X Lacoux, A Morla, G Choquet-Kastylevsky, J Lemoine, Clinical quantitation of prostate-specific antigen biomarker in the low nanogram/milliliter range by conventional bore liquid chromatography-tandem mass spectrometry (multiple reaction monitoring) coupling and correlation with ELISA tests multiple hypothesis testing in microarray experiments. *Mol. Cell Proteomics*. **8**(5), 1006–1015 (2009)
4. C Huillet, A Adrait, D Lebert, G Picard, M Trauchessec, M Louwagie, A Dupuis, L Hittinger, B Ghaleb, P Le Corvoisier, M Jaquinod, J Garin, C Bruley, V Brun, Accurate quantification of cardiovascular biomarkers in serum using protein standard absolute quantification (PSAQ) and selected reaction monitoring. *Mol. Cell Proteomics*. **11**(2) (2012)
5. K Harris, M Girolami, H Mischak, Definition of valid proteomic biomarkers: a Bayesian solution. *Lett. Notes Comput. Sci.* **5780**, 137–149 (2009)
6. M Frantzi, A Bhat, A Latosinska, Clinical proteomic biomarkers: relevant issues on study design & technical considerations in biomarker development. *Clin. Transl. Med.* **3**(7) (2014)
7. P Roy, C Truntzer, D Maucourt-Boulch, T Jouve, N Molinari, Protein mass spectra data analysis for clinical biomarker discovery: a global review. *Brief. Bioinform.* **12**(2), 176–186 (2011)
8. D Sidransky, S Srivastava, Changes in collagen metabolism in prostate cancer: a host response that may alter progression. *Nat. Rev. Cancer*. **18**(3), 789–795 (2003)
9. H Hoijtink, I Klugkist, Comparison of hypothesis testing and Bayesian model selection. *Qual. Quant.* **41**, 73–91 (2007)
10. S Dudoit, J Popper Shaffer, J Boldrick, Multiple hypothesis testing in microarray experiments. *Statist. Sci.* **18**, 71–103 (2003)
11. Y Benjamin, Y Hochberg, Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Statist. Soc. B.* **57**(1), 289–300 (1995)
12. GK Smyth, Linear models and empirical Bayes methods for assessing differential expression in microarray experiments. *Stat. Appl. Genet. Mol. Biol.* **3**(3) (2004)
13. N Draper, H Smith, *Applied regression analysis*, 3rd ed. (Wiley Series in Probability and Statistics, Chichester, New York, Singapore, Toronto, 1998)
14. M Bhattacharjee, C Botting, M Sillanpaa, Bayesian biomarker identification based on marker-expression proteomics data. *ELSEVIER Genomics*. **92**, 37–55 (2008)
15. J Fan, R Li, Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Stat. Assoc.* **96**(456), 1348–1360 (2001)
16. M Chen, Dc Dey, Variable selection for multivariate logistic regression model. *J. Stati. Plan. and Infer.* **111**, 37–55 (2003)
17. Z Yuan, D Ghosh, Combining multiple biomarker models in logistic regression. *Biometrics*. **64** (2008)
18. H Akaike, Information theory and an extension of the maximum likelihood principle. *Proc. Second Int. Symp. Inform.*, 261–281 (1973)
19. AE Hoerl, RW Kennard, Ridge regression: applications to nonorthogonal problems. *Technometrics*. **12**(1), 69–82 (1970)
20. R Tibshirani, Regression shrinkage and selection via the LASSO. *J. Royal Stat. Soc. Series B (Methodology)*. **1**, 267–288 (1996)
21. H Zou, T Hastie, Regularization and variable selection via the elastic net. *J. R. Statist. Soc. B.* **67**, 301–320 (2005)
22. J Ogutu, T Schultz-Streek, HP Piepho, Genomic selection using regularized linear regression models: ridge regression, LASSO, elastic net and their extensions. ser. Proc. of the 15th European workshop on QTL mapping and marker assisted selection (QTLMAS), Rennes, France, 2011, 37–55
23. MI Ghahramani, A note on the evidence and Bayesian Occam's razor. Gatsby Unit, University College London, Technical Report GCNU-TR 2005-003 (2005)
24. CP Robert, G Casella, *Monte-Carlo statistical methods*, ser. Springer Texts in Statistics. (Springer, New York, 2004)
25. B Carlin, T Louis, *Bayesian methods for data analysis*. (CRC Press, Chapman & Hall, Boca Raton, London, New York, 2009)
26. A Raftery, M Newton, J Satagopan, P Krivitsky, in *Bayesian Statistics*. Estimating the integrated likelihood via posterior simulation using the harmonic mean identity, vol. 8, (2007), pp. 1–45
27. D Lee, N Chia, A particle algorithm for sequential Bayesian parameter estimation and model selection. *IEEE Trans. on Sign. Proc.* **50**(2), 326–336 (2002)
28. H Mallick, N Yi, Bayesian methods for high dimensional linear models. *J. Biom. Biostat.* **1**(5), 326–336 (2013)
29. F Adjed, JF Giovannelli, A Giremus, N Dridi, P Szacherski, in *ser. Proc. of the IEEE Intl. Conf. on Acoustics, Speech and Signal Processing (ICASSP 2013)*. Variable selection for a mixed population applied in proteomics (Vancouver, 2013), pp. 1153–1157
30. C Robert, in *The Bayesian Choice. From decision-theoretic foundations to computational implementation*. Springer Text in Statistics (Springer Verlag, New York, 2007)
31. G Saporta, *Probabilités, analyse de données et statistique*, Technip, Ed. Editions TECHNIP, (1990)
32. A Bhattacharyya, On a measure of divergence between two statistical populations defined by their probability distributions. *Indian J. Stat.* **7**(4), 401–406 (1946)
33. L Song, A Smola, A Gretton, J Bedo, K Borgwardt, Feature selection via dependence maximization. *J. Mach. Learn. Res.* **13**(1), 1393–1434 (2012)
34. P Szacherski, JF Giovannelli, L Gerfault, P Mahé, JP Charrier, A Giremus, B Lacroix, P Grangeat, Classification of proteomic MS data as Bayesian solution of an inverse problem. *IEEE Access*. **2**, 1248–1262 (2014)
35. A Klich, C Mercier, L Gerfault, P Grangeat, C Beaulieu, E Degout-Charmette, T Fortin, P Mahé, JF Giovannelli, JP Charrier, A Giremus, D Maucourt-Boulch, P Roy, Experimental design and statistical analysis for evaluation of quantification performance of two molecular profile reconstruction algorithms used in selected reaction monitoring-mass spectrometry. Service de Biostatistique, Hospices Civils and Laboratoire de Biométrie et Biologie Evolutive, Lyon, Technical Report (2016)
36. J Lemoine, T Fortin, A Salvador, A Jaffuel, JP Charrier, G Choquet-Kastylevsky, The current status of clinical proteomics and the use of MRM and MRM3 for biomarker validation. *Pharmacogenomic Proteomic Metabolomic Appl.* **12**(4), 333–345 (2012)

Submit your manuscript to a SpringerOpen® journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► [springeropen.com](http://springeropen.com)