

Multiple Linear Regression: Bayesian Inference for Distributed and Big Data in the Medical Informatics Platform of the Human Brain Project

Lester Melie-Garcia^{1✉}, Bogdan Draganski^{1,2}, John Ashburner³, Ferath Kherif¹

¹ LREN, Department of Clinical Neurosciences, Lausanne University Hospital (CHUV), Switzerland

² Max-Planck-Institute for Human Cognitive and Brain Sciences, Leipzig, Germany

³ Wellcome Trust Centre for Imaging Neuroscience, University College London (UCL), United Kingdom (UK)

✉ Correspondence: Mont Pausible 16, 1011 Lausanne, Switzerland. Phone: +41 21 314 9593. Fax: +41 21 314 1256;
E-mail address: lester.melie@gmail.com

ABSTRACT: We propose a Multiple Linear Regression (MLR) methodology for the analysis of distributed and Big Data in the framework of the Medical Informatics Platform (MIP) of the Human Brain Project (HBP). MLR is a very versatile model, and is considered one of the workhorses for estimating dependences between clinical, neuropsychological and neurophysiological variables in the field of neuroimaging. One of the main concepts behind MIP is to federate data, which is stored locally in geographically distributed sites (hospitals, customized databases, etc.) around the world. We restrain from using a unique federation node for two main reasons: first the maintenance of data privacy, and second the efficiency in management of big volumes of data in terms of latency and storage resources needed in the federation node. Considering these conditions and the distributed nature of data, MLR cannot be estimated in the classical way, which raises the necessity of modifications of the standard algorithms. We use the Bayesian formalism that provides the armamentarium necessary to implement the MLR methodology for distributed Big Data. It allows us to account for the heterogeneity of the possible mechanisms that explain data sets across sites expressed through different models of explanatory variables. This approach enables the integration of highly heterogeneous data coming from different subjects and hospitals across the globe. Additionally, it offers general and sophisticated ways, which are extendable to other statistical models, to suit high-dimensional and distributed multimodal data. This work forms part of a series of papers related to the methodological developments embedded in the MIP.

Keywords: multiple linear regression, general linear model, linear model, linear regression, Bayesian linear model, Bayesian linear regression, MLR, parallel computation, distributed computation, Human Brain Project, Bayesian modeling, model averaging, variational Bayes.

INTRODUCTION

The amount of neuroimaging data (i.e. MRI, fMRI, PET, EEG, etc.) continues to expand exponentially with the production of thousands of studies worldwide as part of clinical and research activities. As an example, after just a few decades of the Magnetic Resonance Imaging technique (MRI), hundreds of thousands of individuals of all ages and conditions have been scanned and the rate of MRI data collection is set to grow considerably over the coming years. Unfortunately, the majority of this data is locally warehoused in laboratories and hospitals and is therefore poorly capitalized into new knowledge, scientific publications and the development of novel medical therapies. This raises the urgent need for a more efficient exploitation of such limitless richness of information in order to shed light on the principles of brain anatomy and function in both healthy and pathological states. In recent years, in light of this necessity, the neuroscientific and software developer's communities have joined efforts to embrace new paradigms related to data sharing, virtualization and federation. These paradigms provide a framework to integrate large amounts of data that will undoubtedly help in generating more realistic models of brain function. Data integration is an effective approach devoted to appropriately combining neuropsychological, genetic, clinical and neuroimaging multimodal data (at the micro, meso and macro scales), in a very heterogeneous population coming from different parts of the globe under dissimilar conditions (race, education level, dietary conditions, etc.).

Precisely, the Medical Informatics Platform (MIP) as part of the Human Brain Project (Subproject 8 (SP8)) (<https://www.humanbrainproject.eu/en/medicine/medical-informatics-platform/>), is conceived to provide an informatics infrastructure to the neuroscientists and clinicians studying brain anatomy and function using clinical along with multimodal big data. This platform faces the challenges of adapting, designing and implementing novel technologies in terms of hardware (clusters etc.), software (i.e. NoSQL, Apache Spark) and data analytics (i.e. decentralized machine learning) to deal with the large and geographically distributed nature of the data. In this paper, we propose a theoretical framework based on the ‘Bayesian formalism’ that contributes to the development of mathematical tools to deal with big and distributed data as part of the MIP environment. Rather than introducing a fully new mathematical development, we gather together and adapt general advances previously described in the Bayesian modeling literature to ultimately propose our methodological framework. We extended these developments to the versatile Multiple Linear Regression (MLR), which is considered one of the workhorses of statistics and data modeling in the neuroimaging field. The Bayesian formalism is ideal for the core of our theoretical machinery since it deals, in a natural way, with the problem of having big, heterogeneous and distributed data.

The paper in general is organized as follows. The main section, ‘Theory’, is divided into two main subsections. The first is devoted to providing the grounds of the general Bayesian framework for distributed data in two ‘regimes’: Parallel and Streaming. A general schema is illustrated that defines the local (hospitals) and federation (global) nodes; the latter federates and delivers the information to the end user. Based on previous works, a model selection/averaging step is proposed, constituting the third level of Bayesian inference. The second subsection describes a special case of the general principles outlined in the former subsection, which are applied to the multiple linear regression model. In addition, a summary of the general equations to be used in practice is shown.

THEORY

1. Problem Statement

One of the main goals of data analysis is to identify the models that best explain the data, along with their intrinsic parameter distributions, in order to discover basic organizational and functional principles of natural systems. Conventional modeling setups and algorithms presuppose that data is stored in one single location as a unique data source (i.e. database). Therefore, new techniques are required to model data when it is large and geographically distributed.

2. General Bayesian computational framework: Parallel approach.

Figure 1 shows the general scheme of our computational framework for distributed data. Data is dispersed across different remote sites, which in our case are hospitals/databases that may be located in different countries. The upper part of the figure represents the Federation node (aka master node, centralized node) that federates the aggregated information (or aggregates) coming from local nodes or hospitals (aka decentralized nodes, worker nodes) represented in the lower part of the figure. Within the framework of this paper, these aggregates are expectations of various model parameters. In this case, the federated estimation of parameter distributions is obtained using local aggregates provided in parallel by all local nodes. This setting is denominated as ‘parallel’ (Broderick et al., 2013).

In this paradigm, the information can be interchanged in both directions. Within an iterative algorithm, the federation node provides intermediate parameter estimates $\hat{\Theta}$ downstream and receives aggregates $(\Theta^{(i)})$ from each local node. The set of aggregates from Hospital ‘ i ’ is represented by $\Theta^{(i)} = [\theta_1^{(i)}, \dots, \theta_N^{(i)}]$ and y_i is the data for the hospital ‘ i ’, where H is the number of hospitals so that $\mathbf{y} = [y_1 \dots, y_H]$.

We proposed the use of the Bayesian formalism as a theoretical framework to define our computational engine, which we named ‘Bayesian computational framework’ (BCF). Therefore our BCF is based on the fundamental Bayes theorem so that the estimation of the parameters $\hat{\Theta}$ at the federation node is expressed by the following equation:

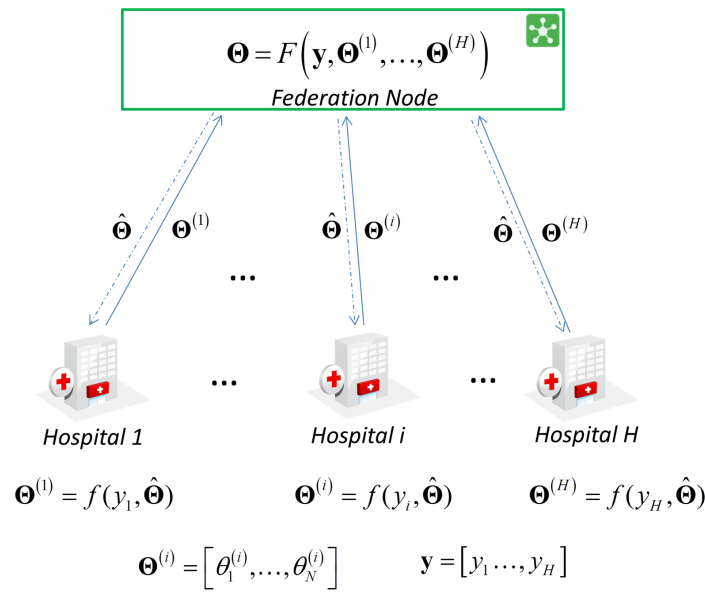


Figure 1. General scheme of the parallel paradigm in the Medical Informatics Platform (Human Brain Project).

$$p(\Theta|y, M_k) \propto p(y|\Theta, M_k) p(\Theta) \quad (1)$$

where $p(\Theta|y, M_k)$ is the posterior distribution of parameters Θ for the specific model M_k and $\hat{\Theta}$ the set of parameter values that maximize $p(\Theta|y, M_k)$; $p(y|\Theta, M_k)$ is the likelihood (taking into account data from all Hospitals) and $p(\Theta)$ the prior probability of the parameters of interest.

Assuming that chunks of data belonging to each hospital are independent, the likelihood can be expressed as:

$$p(y|\Theta, M_k) = p(y_1|\Theta, M_k) p(y_2|\Theta, M_k) \dots p(y_H|\Theta, M_k) \quad (2)$$

Substituting equation (2) into (1), we obtain:

$$p(\Theta|y, M_k) \propto p(y_1|\Theta, M_k) p(y_2|\Theta, M_k) \dots p(y_H|\Theta, M_k) p(\Theta) \quad (3)$$

Multiplying and dividing $M-1$ times by the prior probability $p(\Theta)$ in equation (3) leads to:

$$p(\Theta|y, M_k) \propto \frac{\overbrace{p(y_1|\Theta, M_k) p(\Theta)}^{p(\Theta|y_1, M_k)} \overbrace{p(y_2|\Theta, M_k) p(\Theta)}^{p(\Theta|y_2, M_k)} \dots \overbrace{p(y_H|\Theta, M_k) p(\Theta)}^{p(\Theta|y_H, M_k)}}{(p(\Theta))^{H-1}} \quad (4)$$

$$p(\Theta|y, M_k) \propto \frac{p(\Theta|y_1, M_k) p(\Theta|y_2, M_k) \dots p(\Theta|y_H, M_k)}{(p(\Theta))^{H-1}} = \frac{1}{(p(\Theta))^{H-1}} \prod_{j=1}^H p(\Theta|y_j, M_k) \quad (5)$$

Then the posterior probability of the parameters Θ for a specific model M_k at the federation node is expressed as a product of the posterior probabilities coming from the ‘ H ’ Hospitals. The maximum of the posterior distribution for Θ is achieved at $\hat{\Theta}$ values.

2.1 General Bayesian computational framework: Streaming approach.

In many situations, providing a federated estimate of the parameters of interest when all local estimates are ready to use is very inefficient. Such situations become more critical when there are long delays between local parameter estimations. Hence it has been proposed, when it is tractable, the use of a subset of local estimates to obtain intermediate federated results.

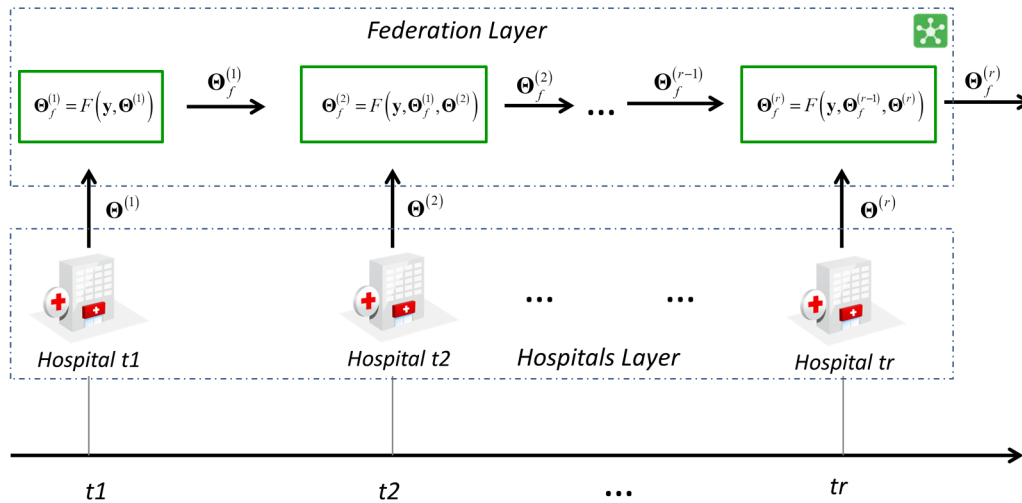


Figure 2. General scheme of the streaming paradigm in the Medical Informatics Platform (Human Brain Project). At the federation layer the new local estimates coming from the hospitals layer are combined with previous ones to deliver federated results.

Therefore as a new local estimate arrives to federation node an update of the federated estimation is performed, without recalculating parameters over again. This general scheme (Figure 2) is known as ‘Streaming’ (Broderick et al., 2013), and under the Bayesian formalism can be expressed through the following equations:

$$\begin{aligned}
 & - p(\Theta | y_1, M_k) \propto p(y_1 | \Theta, M_k) p(\Theta) \\
 & - p(\Theta | y_1, y_2, M_k) \propto p(y_1, y_2 | \Theta, M_k) p(\Theta) = p(y_2 | \Theta, M_k) \overbrace{p(y_1 | \Theta, M_k) p(\Theta)}^{p(\Theta | y_1, M_k)} \\
 & \quad = p(y_2 | \Theta, M_k) p(\Theta | y_1, M_k) \\
 & \quad \vdots \\
 & - p(\Theta | y_1, \dots, y_r, M_k) \propto p(y_1, \dots, y_r | \Theta, M_k) p(\Theta) = \\
 & \quad \quad \quad \underbrace{p(\Theta | y_1, \dots, y_{r-1})}_{p(\Theta | y_1)} \\
 & \quad = p(y_r | \Theta, M_k) p(y_{r-1} | \Theta, M_k) \cdots p(y_1 | \Theta, M_k) p(\Theta) = \\
 & \quad = p(y_r | \Theta, M_k) p(\Theta | y_1, \dots, y_{r-1}, M_k)
 \end{aligned} \tag{6}$$

From the recurrence equations in (6) we have that a new posterior at the federation node, when a specific ‘ r ’ data is ready from a local site, will be the multiplication of the previous posterior distribution by the likelihood function of the new data. Since the data cannot be distributed out of local sites, the posterior

$p(\Theta | \mathbf{y}_1, \dots, \mathbf{y}_{K-1})$ should be sent to the local site ‘ r ’ for updating the federated posterior that is finally, in a new version, pull back to the federation node.

On the other hand, equation (6) can also be rewritten as:

$$p(\Theta | \mathbf{y}, M_k) \propto \frac{\overbrace{p(\mathbf{y}_r | \Theta, M_k) p(\Theta)}^{p(\Theta | \mathbf{y}_r)}}{p(\Theta)} p(\Theta | \mathbf{y}_1, \dots, \mathbf{y}_r, M_k) = \frac{p(\Theta | \mathbf{y}_r, M_k) p(\Theta | \mathbf{y}_1, \dots, \mathbf{y}_r, M_k)}{p(\Theta)} \quad (7)$$

Therefore the new posterior can in addition be computed by equation (7) passing the posterior of parameters Θ from the ‘ r ’th hospital to the federation node to be accordingly combined with the previous posterior. That way is more convenient since the communication between federation and local nodes is preserved in only one direction, which turns out in practice to be simpler to implement.

2.2 Model selection and model averaging: general definitions

As we stated at the beginning of the THEORY section, the selection of one or a set of models that explain our data is key in data analysis. In this subsection we treated the equations for the third level of inference in the Bayesian formalism known as model selection and model averaging (Hoeting et al., 1999; MacKay, 1992; Penny et al., 2006) applied to our particular case.

Given a set of ‘ K ’ candidate models M_1, \dots, M_K for explaining our data, we need to calculate the posterior probability of the models to ultimately either select the best one or apply a model averaging step. In our case we assume the model selection/model averaging to be performed at the federation node, though, this can be implemented at each local site to finally be combined at the federation node. If we assume the data is generated by the same model M_k in all local sites we are in the presence of a Fixed Effects (FFX) analysis. Otherwise if a more relaxed and general assumption is adopted, allowing for the possibility that different data sites use different models, a Random Effect (RFX) analysis is carried out. In the following subsections we separately analyzed both types of analysis.

2.2.1 Model selection: Fixed Effects (FFX) Analysis

The probability distribution for specific model ‘ r ’ is expressed by the well-known equation:

$$p(M_k | \mathbf{y}) = \frac{p(\mathbf{y} | M_k) p(M_k)}{\sum_{i=1}^K p(\mathbf{y} | M_i) p(M_i)} \quad (8)$$

where $p(\mathbf{y} | M_k)$ is the evidence of the model or the marginal likelihood that accounts for all data chunks associated with local sites, $p(M_k)$ the prior probability of M_k that expresses our belief or prior knowledge about the occurrence of the model M_k . The evidence of the model M_k is expressed by:

$$p(\mathbf{y} | M_k) = \int p(\mathbf{y}, \Theta, M_k) d\Theta = \int p(\mathbf{y} | \Theta, M_k) p(\Theta | M_k) d\Theta \quad (9)$$

Because of the condition of independence across data sites expressed in equation (2) we have that $p(M_k | \mathbf{y})$ is:

$$p(M_k | \mathbf{y}) = \frac{p(\mathbf{y} | M_k) p(M_k)}{\sum_{j=1}^K p(\mathbf{y} | M_j) p(M_j)} = \frac{p(\mathbf{y}_1 | M_k) \cdots p(\mathbf{y}_H | M_k) p(M_k)}{\sum_{j=1}^K p(\mathbf{y}_1 | M_j) \cdots p(\mathbf{y}_H | M_j) p(M_j)} \quad (10)$$

$$p(M_k | \mathbf{y}) = \frac{\overbrace{p(\mathbf{y}_1 | M_k) p(M_k)}^{p(M_k | \mathbf{y}_1)} \cdots \overbrace{p(\mathbf{y}_i | M_k) p(M_k)}^{p(M_k | \mathbf{y}_i)} \cdots \overbrace{p(\mathbf{y}_H | M_k) p(M_k)}^{p(M_k | \mathbf{y}_H)}}{\sum_{j=1}^K p(\mathbf{y}_1 | M_j) \cdots p(\mathbf{y}_H | M_j) p(M_j)} \left(\frac{1}{p(M_k)} \right)^{H-1}$$

From equation (10) the posterior probability distribution of a specific model M_k , at the federation node is proportional to the product of the probabilities of this model from all local data sites.

In order to select the model at federation node, that best suits the data across all local sites, we made use of the Group Bayes Factor function as has been proposed by other authors (Stephan et al., 2007):

$$BF_{ij} = \frac{p(M_i | \mathbf{y})}{p(M_j | \mathbf{y})} \quad (11)$$

Assuming uniform model priors $p(M_1) = p(M_2) = \dots = p(M_K) = cte$ equation (11) turns:

$$BF_{ij} = \frac{p(M_i | \mathbf{y})}{p(M_j | \mathbf{y})} = \prod_{h=1}^H \frac{p(M_i | \mathbf{y}_h)}{p(M_j | \mathbf{y}_h)} = \prod_{h=1}^H \frac{p(\mathbf{y}_h | M_i)}{p(\mathbf{y}_h | M_j)} = \prod_{h=1}^H BF_{ij}^{(h)} \quad (12)$$

where $BF_{ij}^{(h)}$ is the Bayes factor for the local data site 'h':

$$BF_{ij}^{(h)} = \frac{p(\mathbf{y}_h | M_i)}{p(\mathbf{y}_h | M_j)} \quad (13)$$

Group Bayes factor encodes the relative probability that the data were generated by one model relative to another assuming that all local data were produced by the same model. It's assumed in practice that for $BF_{ij} > 20$ provides strong evidence in favor of model M_i over M_j .

2.2.2 Model selection: Random Effects (RFX) Analysis

In this subsection we assume the local sites use different models to fit the data that defines a Random Effect (RFX) analysis. RFX accounts for the presence of different underlying mechanisms that explain the data. This approach has been proven (Penny et al., 2010; Stephan et al., 2009) to outperform the classical FFX since its markedly more robustness in the presence of outlier data.

We adopted the methodology developed in Stephan et al. (Stephan et al., 2009) and later extended for model averaging (for family of models) by Penny et. al. (Penny et al., 2010). The proposed RFX approach determines the probability density from which the models that generate the data are sampled at different local sites.

Given the 'K' candidate models, we define a labeling process binary variable or model switch that helps us to calculate the probability that a specific hospital data was generated by the model M_k . Hence, matrix S is defined as:

$$\mathbf{S} = \begin{bmatrix} s_{11} & \cdots & s_{1K} \\ \vdots & \ddots & \vdots \\ s_{H1} & \cdots & s_{HK} \end{bmatrix} \quad (14)$$

where ‘ K ’ is the number of models and ‘ H ’ number of data sites.

$$p(\mathbf{S}) = \prod_{j=1}^H p(\mathbf{s}_j) \quad (15)$$

The row vector $\mathbf{s}_j = [s_{j1}, \dots, s_{jK}]$ expresses the indicator variable to pick one of the ‘ K ’ models for the hospital ‘ j ’. We assume a priori that the distribution of this variable is independent across local sites, and $p(\mathbf{s}_j)$ is in general a Multinomial distribution so that:

$$p(\mathbf{s}_j) = \text{Mult}(\pi_{kj}) \propto \prod_{k=1}^K \pi_{jk}^{s_{jk}} \quad (16)$$

In our case $\sum_{k=1}^K s_{jk} = 1$ that guarantees that only one model will be selected reducing the prior to a Categorical distribution. The element π_{jk} is the prior probability the data in hospital ‘ j ’ was generated by the model ‘ k ’. Therefore the prior probability of \mathbf{S} is finally expressed as:

$$p(\mathbf{S}) = \prod_{j=1}^H p(\mathbf{s}_j) = \prod_{j=1}^H \text{Mult}(\pi_{kj}) \propto \prod_{j=1}^H \prod_{k=1}^K \pi_{jk}^{s_{jk}} \quad (17)$$

The model frequency variable represented by $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$ expresses the frequency of each of the ‘ K ’ models, where $\sum_{k=1}^K \pi_k = 1$. We assume that $p(\boldsymbol{\pi})$ is a Dirichlet distribution so that:

$$p(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\lambda}) \quad (18)$$

with $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_K]$. If a symmetric Dirichlet distribution is assumed, then $\lambda_1 = \lambda_2 = \lambda_3 \dots = \lambda_K = \lambda_0$, and:

$$p(\boldsymbol{\pi}) = \frac{\Gamma(K\lambda_0)}{\Gamma(\lambda_0)^K} \prod_{k=1}^K \pi_k^{\lambda_0-1} \quad (19)$$

The posterior probability $p(\boldsymbol{\pi}|\mathbf{y}, \boldsymbol{\lambda})$ provides the information required to perform model selection or model averaging at the federation node, accounting for model variability across local sites.

We used the Variational Bayes approach (see later in next subsections) to approximate $p(\boldsymbol{\pi}|\mathbf{y}, \boldsymbol{\lambda})$, which is based on the model evidence shown in equation (9). Further details about the derivation of these equations, and an algorithm implementation for estimating the parameters of interest, can be found in (Stephan et al., 2009). Hereafter, only the final equations will be provided.

For model indicator variable \mathbf{s} :

$$q(\mathbf{s}_l) = \text{Mult}(\tilde{\gamma}_m^l) \propto \prod_{m=1}^K (\tilde{\gamma}_m^l)^{s_{lm}} \quad (20)$$

$\tilde{\gamma}_m^l$ is the mean of the multinomial distribution. This is the expected number of times the model ‘ m ’ is responsible of generating data ‘ l ’. The expression for $\tilde{\gamma}_m^l$ is the following:

$$\tilde{\gamma}_m^l = \tilde{\pi}_m \cdot \exp\left[\log\left(p(y_l | s_{lm})\right)\right] \quad (21)$$

where $\log \tilde{\pi}_m = \Psi(\lambda_m) - \Psi\left(\sum_{k=1}^K \lambda_k\right)$, $\Psi(\cdot)$ is the digamma function. Finally, $\tilde{\gamma}_m^l$ is normalized to obtain γ_m^l that is used in the rest of parameter estimators:

$$\gamma_m^l = \frac{\tilde{\gamma}_m^l}{\sum_{k=1}^K \tilde{\gamma}_k^l} \quad (22)$$

For model frequency $\boldsymbol{\pi}$:

The posterior $p(\boldsymbol{\pi}|\mathbf{y}, \lambda)$ is approximated by $q(\boldsymbol{\pi}|\mathbf{y})$, which takes the form of a Dirichlet distribution:

$$q(\boldsymbol{\pi}) = \text{Dir}(\hat{\boldsymbol{\lambda}}) \quad (23)$$

where $\hat{\boldsymbol{\lambda}} = [\hat{\lambda}_1, \dots, \hat{\lambda}_K]$, with $\hat{\lambda}_k$ defined as:

$$\hat{\lambda}_k = \sum_{j=1}^H \gamma_k^j + \lambda_0 \quad (24)$$

Then $\hat{\boldsymbol{\pi}} = [\hat{\pi}_1, \dots, \hat{\pi}_K]$, with:

$$\hat{\pi}_k = E(\pi_k) = \int q(\pi_k) \pi_k d\pi_k = \frac{\hat{\lambda}_k}{\sum_{i=1}^K \hat{\lambda}_i} \quad (25)$$

The model that best explains the overall data is defined as that with highest probability $\hat{\boldsymbol{\pi}}$, according to equation (25). The probability of model 'k', in a RFX analysis at the federation node, is expressed by:

$$p(M_k|\mathbf{y}) = \hat{\pi}_k \quad (26)$$

The most probable model is obtained by:

$$M_i = \max_k [p(M_1|\mathbf{y}), \dots, p(M_k|\mathbf{y}), \dots, p(M_K|\mathbf{y})] \quad (27)$$

Having estimates of the model probabilities accounting for all model evidences and occurrences (at the federation node) allows us to compare two models, rank them or compare families of models. The last use-case is especially important when comparing model families of different natures.

2.2.3 Model averaging

For models of a similar nature, the marginal distribution of the parameters $\boldsymbol{\Theta}$ that accounts for the uncertainty of selecting a specific model at the federation node is given by:

$$p(\boldsymbol{\Theta}|\mathbf{y}) = \sum_{k=1}^K p(\boldsymbol{\Theta}|\mathbf{y}, M_k) p(M_k|\mathbf{y}) \quad (28)$$

The posterior probability of model k at the federation node $p(M_k|\mathbf{y})$ may be calculated using a FFX - equation (10) - or RFX - equation (26) - analysis. Posterior probabilities of the parameters $p(\boldsymbol{\Theta}|\mathbf{y}, M_k)$ for a specific model would be provided either by the equation (5) or (7).

3. Multiple Linear regression (MLR) for distributed data: probabilistic generative model.

In this subsection we apply the previous general approach to the special case of the multiple linear regression model (MLR). It allows comparing different MLRs as considering subsets of explanatory variables that best explain the dependent data. This approach accounts for distributed data where only aggregates can be passed to the federation node to ultimately estimate the distribution of the parameters that best explained the data across local data sites.

Multiple linear regression allows to study the relationship between a dependent variable \mathbf{y} and the explanatory variables (or independent variables) denoted by \mathbf{A} . The general problem is to solve the standard linear equation posed as following:

$$\mathbf{y} = \mathbf{A}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (29)$$

\mathbf{y} : vector of dependent variable $p \times 1$; $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_N]$: matrix of regressors or design matrix so that $\mathbf{A}_i = [a_{i1}, \dots, a_{iH}]^T$, $p \times N$; $\boldsymbol{\beta}$: vector of regression coefficients $N \times 1$; $\boldsymbol{\varepsilon}$: vector of experimental noise/error, $p \times 1$.

We have the special situation in the MIP that data \mathbf{y} and \mathbf{A} are naturally split and distributed over H different sites (hospitals). Additionally the data cannot be fetched to a federation site to solve the equation (1). Thus we have to adapt classical methods for estimating $\boldsymbol{\beta}$, which requires combining aggregates from hospitals data.

We propose to use the Variational Bayes approach (VB) in order to solve MLR equation (29) at each local site. VB maximizes the negative 'Free energy' (NFE) as a surrogate measure of the log evidence or marginal log likelihood. Using either parallel or streaming paradigms, the MLR model parameters distributions and model evidences are properly combined in the federation node to provide global model parameters distributions, which account for the heterogeneity of data across local data sites.

3.1 MLR equations for local sites (Hospitals): Hierarchical Bayesian Model

The MLR equation (29) is solved at each local data site. Therefore we omitted the sub index referring to a specific data site. The equations derived in this subsection are applied in the same way for each local site. In the following a VB formalism will be developed for MLR, partially treated by previous authors (Tzikas et al., 2008). The variables' dependences of our MLR formalism is represented via the Directed Acyclic graph (DAG) shown in Figure 3.

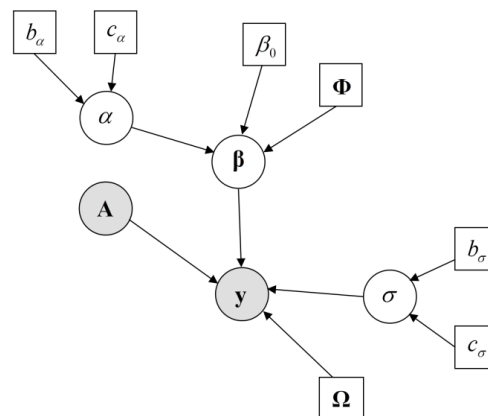


Figure 3. Directed Acyclic graph (DAG) of the MLR. The variables shown in squares are defined a priori. Shaded circles represent observed data whereas not shaded circles represent the model parameters (random variables) to be estimated.

The joint distribution of the data and model parameters for a specific model M_k is given by:

$$p(\mathbf{y}|\boldsymbol{\beta}, \sigma, \alpha|M_k) = p(\mathbf{y}|\boldsymbol{\beta}, \sigma, \boldsymbol{\Omega}, M_k) p(\boldsymbol{\beta}|\alpha, \boldsymbol{\Phi}, M_k) p(\alpha) p(\sigma) \quad (30)$$

Where the vector of model parameters is $\boldsymbol{\Theta} = [\boldsymbol{\beta}, \sigma, \alpha]$; so that σ and α are the data and linear coefficients $\boldsymbol{\beta}$ precisions respectively.

Table 1. Example of the MLR definition for three explanatory variables case. The ones in the table indicate that specific variables are included in the model. In this case we have $N=3$ variables with $K=7$ possible models.

	β_1	β_2	β_3	MLR Equation
M_1	0	0	1	$y = A_3 \beta_3$
M_2	0	1	0	$y = A_2 \beta_2$
M_3	0	1	1	$y = [A_2 \ A_3] \begin{bmatrix} \beta_2 \\ \beta_3 \end{bmatrix}$
M_4	1	0	0	$y = A_1 \beta_1$
M_5	1	0	1	$y = [A_1 \ A_3] \begin{bmatrix} \beta_1 \\ \beta_3 \end{bmatrix}$
M_6	1	1	0	$y = [A_1 \ A_2] \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$
M_7	1	1	1	$y = [A_1 \ A_2 \ A_3] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$

A specific model M_k is defined as a subset of explanatory variables. Among models the full set of variables is included whereas the absence of any explanatory variable (null model) is excluded. Thus the maximum number of models is $K = 2^N - 1$. Hence, matrix \mathbf{A} changes the number of columns for each model. Table 1 shows an example of $N=3$ explanatory variables and the seven possible models.

3.2 Likelihood

We assumed \mathbf{y} is normally distributed as follows:

$$p(\mathbf{y}|\boldsymbol{\beta}, \sigma, \boldsymbol{\Omega}, M_k) = N(\mathbf{A} \boldsymbol{\beta}, \sigma^{-1} \boldsymbol{\Omega}) = \frac{\sigma^{p/2}}{(2\pi)^{p/2} |\boldsymbol{\Omega}|^{1/2}} e^{-\frac{\sigma}{2} (\mathbf{y} - \mathbf{A} \boldsymbol{\beta})^T \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{A} \boldsymbol{\beta})} \quad (31)$$

The experimental error (noise) is considered to have a normal distribution so that:

$$p(\boldsymbol{\epsilon}) = N(\mathbf{0}, \sigma^{-1} \boldsymbol{\Omega}) \quad (32)$$

where σ is the error precision; $\boldsymbol{\Omega}$ is the covariance matrix modeling the correlation structure present in our data.

3.3 Prior distributions

This subsection provides the mathematical form of the prior probability distributions assumed for the parameters of interest.

Regression coefficients β :

$$p(\beta|\alpha, \Phi) = N(\beta_0, \alpha^{-1}\Phi) = \frac{\alpha^{N/2}}{(2\pi)^{N/2} |\Phi|^{1/2}} e^{-\frac{\alpha}{2}(\beta - \beta_0)^T \Phi^{-1} (\beta - \beta_0)} \quad (33)$$

Precision of the regression coefficients α :

$$p(\alpha) = Ga(b_\alpha, c_\alpha) = \frac{c_\alpha^{b_\alpha}}{\Gamma(b_\alpha)} \alpha^{b_\alpha-1} e^{-c_\alpha \alpha} \quad (34)$$

Precision of the likelihood function σ :

$$p(\sigma) = Ga(b_\sigma, c_\sigma) = \frac{c_\sigma^{b_\sigma}}{\Gamma(b_\sigma)} \sigma^{b_\sigma-1} e^{-c_\sigma \sigma} \quad (35)$$

3.4 Posteriors

Our aim is to estimate the posterior distributions of the unknown parameters $\Theta = [\beta, \sigma, \alpha]$ and the evidence of the different models M_1, \dots, M_K . In doing so, we use the *Variational Bayes* formalism (VB) (Beal, 2003; Lappalainen and Miskin, 2000). This is an efficient method already used in previous publications (Penny et al., 2005; Trujillo-Barreto et al., 2008). It is based on the Variational Free Energy method of Feynman and Bogoliubov, which provides an approximate factorized posterior distribution $q(\Theta|\mathbf{Y})$, with minimal Kullback–Leibler (KL) divergence from the true posterior $p(\Theta|\mathbf{y}) = p(\beta, \sigma, \alpha|\mathbf{y})$.

We consider the following factorization of the approximate posterior distribution:

$$q(\beta, \sigma, \alpha|\mathbf{y}) = q(\beta|\mathbf{y}) q(\sigma|\mathbf{y}) q(\alpha|\mathbf{y}) \quad (36)$$

The following subsections provide details about the derivation of the approximate posteriors for each parameter. In general, the approximate posterior $q(\Theta|\mathbf{Y})$ for the i -th subset of parameters is expressed by $q(\theta_i|\mathbf{Y})$ as:

$$\begin{aligned} q(\theta_i|\mathbf{Y}) &= \frac{\exp\left(\int q(\Theta_{\setminus i}|\mathbf{Y}) \log p(\mathbf{Y}, \Theta) d\Theta_{\setminus i}\right)}{\int \exp\left(\int q(\Theta_{\setminus i}|\mathbf{Y}) \log p(\mathbf{Y}, \Theta) d\Theta_{\setminus i}\right) d\theta_i} \\ &= \frac{\exp(I(\theta_i))}{\int \exp(I(\theta_i)) d\theta_i} \end{aligned} \quad (37)$$

where $I(\theta_i) = \int q(\Theta_{\setminus i}|\mathbf{Y}) \log p(\mathbf{Y}, \Theta) d\Theta_{\setminus i} = \langle \log p(\mathbf{Y}, \Theta) \rangle_{q(\Theta_{\setminus i})}$, and $\Theta_{\setminus i}$ indicates parameters not present in the i -th group of parameters.

3.4.1 Posterior for regression coefficients β

The approximated posterior distribution $q(\beta|y)$ of the regression coefficients β is expressed by:

$$q(\beta|y) \propto e^{I(\beta)} \quad (38)$$

where

$$I(\beta) = \int q(\sigma|y) q(\alpha|y) \log [p(y|\beta, \sigma) p(\beta|\alpha) p(\sigma) p(\alpha)] d\sigma d\alpha \quad (39)$$

Finally $q(\beta|y)$ is a normal distribution so that,

$$q(\beta|y) = N(\hat{\beta}, \hat{\Phi}_\beta) = \frac{1}{(2\pi)^{N/2} |\hat{\Phi}_\beta|^{1/2}} e^{-\frac{1}{2}(\beta - \hat{\beta})^T \hat{\Phi}_\beta^{-1} (\beta - \hat{\beta})} \quad (40)$$

The mean and covariance matrix of this multivariate normal distribution is given by:

$$\hat{\beta} = \hat{\Phi}_\beta^{-1} (\hat{\sigma}^T \mathbf{A}^{-1} \mathbf{y} + \hat{\alpha} \Phi^{-1} \beta_0) \quad (41)$$

$$\hat{\Phi}_\beta = \hat{\sigma}^T \mathbf{A}^{-1} \mathbf{A} + \hat{\alpha} \Phi \quad (42)$$

3.4.2 Posterior for the precision of the regression coefficients α

In this case the approximated posterior distribution for α is obtained by:

$$q(\alpha|y) \propto e^{I(\alpha)} \quad (43)$$

where

$$I(\alpha) = \int q(\beta|y) q(\sigma|y) \log [p(y|\beta, \sigma) p(\beta|\alpha) p(\sigma) p(\alpha)] d\beta d\sigma \quad (44)$$

After some algebra solving equation (44) we see that the posterior distribution of α is approximated by a gamma distribution so that:

$$q(\alpha|y) = Ga(\hat{b}_\alpha, \hat{c}_\alpha) \quad (45)$$

where

$$\begin{aligned} \hat{b}_\alpha &= b_\alpha + \frac{N}{2} - 1 \\ \hat{c}_\alpha &= c_\alpha + \frac{1}{2} (\hat{\beta} - \beta_0)^T \Phi^{-1} (\hat{\beta} - \beta_0) + Tr(\Phi^{-1} \hat{\Phi}_\beta) \end{aligned} \quad (46)$$

Finally, the expected value of α is:

$$\hat{\alpha} = \frac{\hat{b}_\alpha}{\hat{c}_\alpha} \quad (47)$$

3.4.3 Posterior for the precision σ

The approximate posterior distribution for σ is obtained by:

$$q(\sigma|\mathbf{y}) \propto e^{I(\sigma)} \quad (48)$$

where

$$I(\sigma) = \int q(\boldsymbol{\beta}|\mathbf{y}) q(\boldsymbol{\alpha}|\mathbf{y}) \log [p(\mathbf{y}|\boldsymbol{\beta}, \sigma) p(\boldsymbol{\beta}|\boldsymbol{\alpha}) p(\sigma) p(\boldsymbol{\alpha})] d\boldsymbol{\beta} d\boldsymbol{\alpha} \quad (49)$$

After some algebra, we obtain the following gamma distribution:

$$q(\sigma|\mathbf{y}) = Ga(\hat{b}_\sigma, \hat{c}_\sigma) \quad (50)$$

where

$$\begin{aligned} \hat{b}_\sigma &= b_\sigma + \frac{p}{2} - 1 \\ \hat{c}_\sigma &= c_\sigma + \frac{1}{2} (\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\beta}})^T \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\beta}}) + Tr(\mathbf{A}^T \boldsymbol{\Omega}^{-1} \mathbf{A} \hat{\boldsymbol{\Phi}}_\beta) \end{aligned} \quad (51)$$

The expected value of σ is given by:

$$\hat{\sigma} = \frac{\hat{b}_\sigma}{\hat{c}_\sigma} \quad (52)$$

3.5 Free Energy and KL divergences

The free energy function is expressed as:

$$F = \int q(\boldsymbol{\Theta}|\mathbf{y}) \log p(\mathbf{y}|\boldsymbol{\Theta}) d\boldsymbol{\Theta} - \int q(\boldsymbol{\Theta}|\mathbf{y}) \log \frac{q(\boldsymbol{\Theta}|\mathbf{y})}{p(\boldsymbol{\Theta})} d\boldsymbol{\Theta} \quad (53)$$

where the first term $\langle \log p(\mathbf{y}|\boldsymbol{\Theta}) \rangle_{q(\boldsymbol{\Theta}|\mathbf{y})} = \int q(\boldsymbol{\Theta}|\mathbf{y}) \log p(\mathbf{y}|\boldsymbol{\Theta}) d\boldsymbol{\Theta}$ is the mean likelihood over $q(\boldsymbol{\Theta}|\mathbf{y})$ and $KL(q(\boldsymbol{\Theta}|\mathbf{y})||p(\boldsymbol{\Theta})) = \int q(\boldsymbol{\Theta}|\mathbf{y}) \log \frac{q(\boldsymbol{\Theta}|\mathbf{y})}{p(\boldsymbol{\Theta})} d\boldsymbol{\Theta}$ the Kullback Leibler divergence between the approximated posterior and prior distributions.

After some algebraic work, the expectation of the Likelihood is given by:

$$\begin{aligned} \langle \log p(\mathbf{y}|\boldsymbol{\Theta}) \rangle_{q(\boldsymbol{\Theta}|\mathbf{y})} &= \frac{N}{2} \left(\Psi(\hat{b}_\sigma) - \ln \hat{c}_\sigma \right) - \frac{N}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Omega}| - \\ &\quad - \frac{\hat{\sigma}}{2} (\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\beta}})^T \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\beta}}) - \frac{\hat{\sigma}}{2} Tr(\mathbf{A}^T \boldsymbol{\Omega}^{-1} \mathbf{A} \hat{\boldsymbol{\Phi}}_\beta) \end{aligned} \quad (54)$$

where $Tr(\cdot)$ is the trace operator.

Similarly, the Kullback Leibler divergence between the approximated posterior and prior distributions is expressed as:

$$\begin{aligned} KL(q(\Theta|y)||p(\Theta)) &= \int q(\beta|y)q(\sigma|y)q(\alpha|y) \log \frac{q(\beta|y)q(\sigma|y)q(\alpha|y)}{p(\beta|\alpha)p(\sigma)p(\alpha)} \\ &= KL(q(\beta|y)||p(\beta|\alpha)) + KL(q(\sigma|y)||p(\sigma)) + KL(q(\alpha|y)||p(\alpha)) \end{aligned} \quad (55)$$

After some algebraic operations we obtain:

$$\begin{aligned} KL(q(\beta|y)||p(\beta|\alpha)) &= -\frac{1}{2} \log \frac{|\hat{\Phi}_\beta|}{|\Phi|} - \frac{N}{2} - \frac{N}{2} (\Psi(\hat{b}_\alpha) - \log \hat{c}_\alpha) + \\ &\quad + \frac{\hat{\alpha}}{2} Tr(\Phi^{-1} \hat{\Phi}_\beta) + \frac{\hat{\alpha}}{2} (\hat{\beta} - \beta_0)^T \Phi^{-1} (\hat{\beta} - \beta_0) \end{aligned} \quad (56)$$

$$\begin{aligned} KL(q(\sigma|y)||p(\sigma)) &= (\hat{b}_\sigma - 1) \Psi(\hat{b}_\sigma) + \log \hat{c}_\sigma - \hat{b}_\sigma - \log \Gamma(\hat{b}_\sigma) + \log \Gamma(b_\sigma) - \\ &\quad - b_\sigma \log(c_\sigma) - (b_\sigma - 1) (\Psi(\hat{b}_\sigma) - \log \hat{c}_\sigma) + \frac{\hat{b}_\sigma c_\sigma}{\hat{c}_\sigma} \end{aligned} \quad (57)$$

$$\begin{aligned} KL(q(\alpha|y)||p(\alpha)) &= (\hat{b}_\alpha - 1) \Psi(\hat{b}_\alpha) + \log \hat{c}_\alpha - \hat{b}_\alpha - \log \Gamma(\hat{b}_\alpha) + \log \Gamma(b_\alpha) - \\ &\quad - b_\alpha \log(c_\alpha) - (b_\alpha - 1) (\Psi(\hat{b}_\alpha) - \log \hat{c}_\alpha) + \frac{\hat{b}_\alpha c_\alpha}{\hat{c}_\alpha} \end{aligned} \quad (58)$$

Equations (54), (56), (57) and (58) provides the necessary quantities to evaluate the Free energy function. The dependence among the posteriors of the model parameters $\Theta = [\beta, \sigma, \alpha]$ in equations (41), (46) and (51) leads to an iterative algorithm. The Free Energy is evaluated at each iteration, and is expected to continue to increase until the convergence criterion is met. This is defined as $\frac{F_i - F_{i-1}}{F_{i-1}} < Tol$, where F_i and F_{i-1} are the free energy at iteration 'i' and 'i-1' respectively, Tol is the tolerance with values defined in $0 \leq Tol \leq 1$.

Since the Free Energy is a surrogate measure of the log model evidence, the model probabilities - equation (8) - at local nodes, are given by the expression:

$$p(M_k|y) = \frac{e^{F(y|M_k)} p(M_k)}{\sum_{i=1}^K e^{F(y|M_i)} p(M_i)} \quad (59)$$

where the model evidence is expressed by:

$$p(y|M_k) \approx e^{F(y|M_k)} \quad (60)$$

As a final step, the model evidences and the distributions of the parameters $\Theta = [\beta, \sigma, \alpha]$ from each local node are submitted to the federation node to be properly federated.

The general iterative algorithm that runs at each local node is the following:

Local Node –General Algorithm -

- Local data \mathbf{y} and \mathbf{A} assessment.
- Definition of the stopping criteria tolerance ' Tol '
- Definition of $\mathbf{\Omega}$, $\mathbf{\Phi}$ as the covariance structure for \mathbf{y} and $\mathbf{\beta}$ respectively.
- Initialize hyperparameters: $b_\sigma, c_\sigma, b_\alpha, c_\alpha$
- For $k=1$ to K (Number of Models)
 - Iteration ' i '
 - Compute $\hat{\alpha}$ using equation (47)
 - Compute $\hat{\sigma}$ using equation (52)
 - Compute $\hat{\mathbf{\beta}}$ and $\hat{\mathbf{\Phi}}_\beta$ using equations (41) and (42) respectively
 - Compute Free Energy F_i for iteration ' i ' combining equations (54), (56), (57) and (58)
 - Repeat until $\frac{F_i - F_{i-1}}{F_{i-1}} < Tol$
 - Evaluate evidence of model M_k using equation (60)
- End (Number of Models)
- Send to the Federation node: evidences of the ' K ' models and distributions of the parameters $\hat{\Theta} = [\hat{\mathbf{\beta}}, \hat{\sigma}, \hat{\alpha}]$

3.6 MLR equations for the federation node

At the federation node, the *MLR* is estimated accounting for parameter distributions and model evidences coming from local nodes. In order to obtain expressions (see next subsection) for the federated *MLR* model parameters we made use of the general equations developed in all subsections 2.x.

3.6.1 Linear model coefficients $\mathbf{\beta}$:

At each local node, this parameter has a normal distribution given by equations (40), (41) and (42). The distribution of $\mathbf{\beta}$ at the federation node using equation (5) is given by:

$$p(\mathbf{\beta}|\mathbf{y}, M_k) = \frac{1}{[p(\mathbf{\beta}|\alpha)]^{H-1}} \prod_{j=1}^H q(\mathbf{\beta}|\mathbf{y}_j, M_k) \quad (61)$$

$$p(\mathbf{\beta}|\mathbf{y}, M_r) = N(\hat{\mathbf{\beta}}_f, \hat{\mathbf{\Phi}}_f) = \frac{1}{(2\pi)^{N/2} |\hat{\mathbf{\Phi}}_f|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{\beta} - \hat{\mathbf{\beta}}_f)^T \hat{\mathbf{\Phi}}_f^{-1} (\mathbf{\beta} - \hat{\mathbf{\beta}}_f) \right] \quad (62)$$

where,

$$\hat{\mathbf{\beta}}_f = \hat{\mathbf{\Phi}}_f^{-1} \left(\sum_{i=1}^H \mathbf{\Phi}_{\beta_i}^{-1} \hat{\mathbf{\beta}}_i - \hat{\alpha}_f (H-1) \mathbf{\Phi}^{-1} \mathbf{\beta}_0 \right) \quad (63)$$

$$\hat{\mathbf{\Phi}}_f = \sum_{i=1}^M \mathbf{\Phi}_{\beta_i}^{-1} - \hat{\alpha}_f (M-1) \mathbf{\Phi}^{-1} \quad (64)$$

3.6.2 Precision of the regression coefficients α

At each local node, α follows a Gamma distribution given by equations (45), (46) and (47). After some algebraic operations on equation (5) we see that α also distributes as gamma at the federation node so that:

$$p(\alpha | \mathbf{y}, M_k) = \frac{1}{[p(\alpha)]^{H-1}} \prod_{j=1}^H q(\alpha | \mathbf{y}_j, M_k) \quad (65)$$

$$p(\alpha | \mathbf{y}, M_k) = Ga(b_{\alpha_f}, c_{\alpha_f}) = \frac{c_{\alpha_f}^{b_{\alpha_f}}}{\Gamma(b_{\alpha_f})} \alpha_f^{b_{\alpha_f}-1} e^{-c_{\alpha_f} \alpha_f} \quad (66)$$

$$\begin{aligned} b_{\alpha_f} &= \sum_{j=1}^H \hat{b}_{\alpha_j} - (H-1)b_{\alpha} \\ c_{\alpha_f} &= \sum_{j=1}^H \hat{c}_{\alpha_j} - (H-1)c_{\alpha} \end{aligned} \quad (67)$$

The expected value for α is expressed by:

$$\hat{\alpha}_f = \frac{b_{\alpha_f}}{c_{\alpha_f}} \quad (68)$$

3.6.3 Precision of the likelihood function σ

At each local node, σ has a Gamma distribution given by equations (50), (51) and (52). After combining the local estimates through equation (5) this model parameter distributes as gamma at the federation node. This probability distribution is expressed as:

$$p(\sigma | \mathbf{y}, M_k) = \frac{1}{[p(\sigma)]^{H-1}} \prod_{j=1}^H q(\sigma | \mathbf{y}_j, M_k) \quad (69)$$

$$p(\sigma | \mathbf{y}, M_k) = Ga(b_{\sigma_f}, c_{\sigma_f}) = \frac{c_{\sigma_f}^{b_{\sigma_f}}}{\Gamma(b_{\sigma_f})} \sigma^{b_{\sigma_f}-1} e^{-c_{\sigma_f} \sigma} \quad (70)$$

$$\begin{aligned} b_{\sigma_f} &= \sum_{j=1}^H \hat{b}_{\sigma_j} - (H-1)b_{\sigma} \\ c_{\sigma_f} &= \sum_{j=1}^H \hat{c}_{\sigma_j} - (H-1)c_{\sigma} \end{aligned} \quad (71)$$

with the expected value given by:

$$\hat{\sigma}_f = \frac{b_{\sigma_f}}{c_{\sigma_f}} \quad (72)$$

3.6.4 Model selection and averaging for MLR at the federation node

As we mentioned before, a specific model M_k in the MLR is defined as a subset of explanatory variables (see Table 1). This means we are looking for a subset of β_i in β that best contributes explaining the data across hospital sites. The number of possible models ' K ' is $K = 2^N - 1$, taking into account that each

explanatory variable can be present or not. The model selection at the federation node is conducted by using either a FFX or RFX approach given by equation (10) or (26) respectively. The model averaging step, via expression (28), provides the estimation of the parameter distributions irrespective of the model selected. For the regression coefficients this equation takes the mathematical form:

$$p(\beta|y) = \sum_{k=1}^K p(\beta|y, M_k) p(M_k|y) \quad (73)$$

The general algorithm that runs at the federation node is the following:

Federation Node –General Algorithm -

- Fetch from the ‘H’ local nodes the local parameter expectations: $\hat{\beta}_j$, Φ_{β_j} , \hat{b}_{α_j} , \hat{c}_{α_j} , \hat{b}_{σ_j} , \hat{c}_{σ_j} for all models $M_k, k=1, \dots, K$.
- Estimation of the federated expected value $\hat{\alpha}_f$ of the precision of the linear model coefficients using equation (68)
- Estimation of the federated expected value $\hat{\sigma}_f$ of the precision of the data distribution using equation (72).
- Estimation of the federated linear coefficients $\hat{\beta}_f$ and the covariance matrix $\hat{\Phi}_f$ using equations (63) and (64) respectively.
- Execution of the Model selection using a ‘Random Effects (RFX) Analysis’ (see the corresponding section in the text).
- Model averaging using equation (73).
- Submit to the user the distributions and its expectations $\hat{\beta}_f$, $\hat{\Phi}_f$, $\hat{\sigma}_f$, $\hat{\alpha}_f$ for the model parameters and the model that best explains the data (set of explanatory variables that best explains the data).

CONCLUDING REMARKS

In this paper we presented a Bayesian formalism to deal with distributed, highly heterogeneous and big data, based on previous theoretical developments. This formalism is part of the Bayesian computational core embedded in the Medical Informatics Platform (MIP) of the European Human Brain Project (www.humanbrainproject.eu). As an application example we applied the general approach outlined in this work to the multiple linear regression (MLR), which is considered one of the workhorses of statistics in the neuroimaging field. Though the adaptation of the MLR to a distributed environment might be considered straightforward, in this work we extended it with a model selection and averaging approach. This is crucial for selecting appropriate models to explain highly heterogeneous data.

The federation node is defined as the site in which the local aggregates quantities are combined to deliver a final result accounting for the diversity of the different hospital data (local nodes). We provided approximated probability distributions of the MLR parameters using the Variational Bayes formalism (VB). VB approximates the parameters probability distributions based on the minimization of the negative free energy (NFE) as a surrogate measure of the log evidence. NFE enables the estimation of model probabilities to ultimately performing the model selection and averaging step at the federation node. In this respect, we brought into our formalism the methodology proposed by Stephan et al. (Stephan et al., 2009) (for ‘Dynamic Causal Modeling’-DCM-) to account for the heterogeneity of the possible mechanisms that explain data at different hospitals. This is an important aspect in order to integrate highly heterogeneous data coming from different subjects and hospitals across the globe. For the case of a large number of explanatory variables ($N \gg 1$) an exhaustive evaluation of all model evidences is computationally prohibitive. In such cases a Markov Chain Monte Carlo (MCMC) approach ought to be adopted. MCMC would sample the mass of the models probability distribution allowing, as before, model selection and model averaging steps to be conducted. In addition, a factorization over β coefficients in the VB scheme -equation (36)- should be assumed at local nodes in order to get tractable estimations.

The implementation of our theoretical formalism and its application, in both simulated and real data, will be addressed in a separate work. Its performance under different conditions and the limits of its application will be evaluated through simulated data. Several situations should be considered: a) unbalanced data across local nodes; b) influence of the number of local nodes on model selection and parameters estimation; c) robustness to the presence of outliers; d) robustness of the estimations when the number of explanatory variables increases, etc.

The present work is the first of a series of papers adapting, through Bayesian formalism, the existing statistical models to distributed and big data in order to ultimately discover underlying mechanisms of brain anatomy and function (using the Medical Informatics Platform). This approach is not limited to the Neuroscience field; it can be applied in other branches of science and technology.

ACKNOWLEDGEMENTS

The research leading to these results has received funding from the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement No 604102 and the European Union's Horizon 2020 research and innovation programme under grant agreement No 720270 (HBP SGA1). BD is supported by the Swiss National Science Foundation (NCCR Synapsy, project grant Nr 32003B_159780 and SPUM 33CM30_140332/1), Foundation Parkinson Switzerland and Foundation Synapsis.

REFERENCES

- Beal, M.J., 2003. Variational algorithms for approximate Bayesian inference. PhD thesis, Gatsby Computational Neuroscience Unit, University College London, London.
- Broderick, T., Boyd, N., Wibisono, A., Wilson, A.C., Jordan, M.I., 2013. Streaming Variational Bayes, in: *Advances in Neural Information Processing Systems*. pp. 1727–1735.
- Hoeting, J.A., Madigan, D., Raftery, A.E., Volinsky, C.T., 1999. Bayesian model averaging: a tutorial. *Stat. Sci.* 382–401.
- Lappalainen, H., Miskin, J., 2000. Ensemble learning, *Advances in Independent Component Analysis*.
- MacKay, D.J.C., 1992. Bayesian Interpolation. *Neural Comput.* 4, 415–447. doi:10.1162/neco.1992.4.3.415
- Penny, W., Mattout, J., Trujillo-Barreto, N., 2006. Bayesian model selection and averaging. *Stat. Parametr. Mapp. Anal. Funct. brain images*. London Elsevier.
- Penny, W.D., Stephan, K.E., Daunizeau, J., Rosa, M.J., Friston, K.J., Schofield, T.M., Leff, A.P., 2010. Comparing Families of Dynamic Causal Models. *PLoS Comput. Biol.* 6, e1000709. doi:10.1371/journal.pcbi.1000709
- Penny, W.D., Trujillo-Barreto, N.J., Friston, K.J., 2005. Bayesian fMRI time series analysis with spatial priors. *Neuroimage* 24, 350–362. doi:10.1016/j.neuroimage.2004.08.034
- Stephan, K.E., Marshall, J.C., Penny, W.D., Friston, K.J., Fink, G.R., 2007. Interhemispheric Integration of Visual Processing during Task-Driven Lateralization. *J. Neurosci.* 27.
- Stephan, K.E., Penny, W.D., Daunizeau, J., Moran, R.J., Friston, K.J., 2009. Bayesian model selection for group studies. *Neuroimage* 46, 1004–1017. doi:10.1016/j.neuroimage.2009.03.025
- Trujillo-Barreto, N.J., Aubert-Vázquez, E., Penny, W.D., 2008. Bayesian M/EEG source reconstruction with spatio-temporal priors. *Neuroimage* 39, 318–335. doi:10.1016/j.neuroimage.2007.07.062
- Tzikas, D., Likas, A., Galatsanos, N., 2008. The variational approximation for Bayesian inference. *IEEE Signal Process. Mag.* 25, 131–146. doi:10.1109/MSP.2008.929620