

Springer Proceedings in Mathematics & Statistics

Antonio Canale
Daniele Durante
Lucia Paci
Bruno Scarpa *Editors*

Studies in Neural Data Science

StartUp Research 2017, Siena, Italy,
June 25–27

 Springer

Springer Proceedings in Mathematics & Statistics

Volume 257

Springer Proceedings in Mathematics & Statistics

This book series features volumes composed of selected contributions from workshops and conferences in all areas of current research in mathematics and statistics, including operation research and optimization. In addition to an overall evaluation of the interest, scientific quality, and timeliness of each proposal at the hands of the publisher, individual contributions are all refereed to the high quality standards of leading journals in the field. Thus, this series provides the research community with well-edited, authoritative reports on developments in the most exciting areas of mathematical and statistical research today.

More information about this series at <http://www.springer.com/series/10533>

Antonio Canale · Daniele Durante
Lucia Paci · Bruno Scarpa
Editors

Studies in Neural Data Science

StartUp Research 2017, Siena, Italy,
June 25–27

 Springer

Editors

Antonio Canale
Department of Statistical Sciences
University of Padova
Padua, Italy

Lucia Paci
Department of Statistical Sciences
Università Cattolica del Sacro Cuore
Milan, Italy

Daniele Durante
Department of Decision Sciences
Bocconi University
Milan, Italy

Bruno Scarpa
Department of Statistical Sciences
University of Padova
Padua, Italy

ISSN 2194-1009 ISSN 2194-1017 (electronic)
Springer Proceedings in Mathematics & Statistics
ISBN 978-3-030-00038-7 ISBN 978-3-030-00039-4 (eBook)
<https://doi.org/10.1007/978-3-030-00039-4>

Library of Congress Control Number: 2018961221

Mathematics Subject Classification (2010): S11001, S17030, B18006, L15020

© Springer Nature Switzerland AG 2018

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Preface

This volume contains a collection of peer-reviewed articles arising from *StartUp Research*. This meeting took place on June 25–27, 2017, at the ancient Certosa di Pontignano (Pontignano Charterhouse), a few kilometers from Siena (Italy). *StartUp Research* was a satellite event of the Statistical Conference of the Italian Statistical Society, held in Florence (Italy) in June 2017. The event was additionally endorsed by the young group of the Italian Statistical Society (<https://youngsis.github.io/>), whose aim is to promote activities and provide a social networking platform for early-career researchers in statistics.

StartUp Research was a stimulating experience. It brought together 28 early-career researchers in statistics and seven international professors with the common task of developing novel statistical methods for complex and multimodal brain imaging data. It is, in fact, increasingly common in neuroscience to monitor the brain activity of each subject under different imaging technologies. This motivates the development of novel statistical methods for joint modeling of complex and multimodality data on brain function and structure. The junior researchers, divided into seven groups, focused on brain imaging data from a study of the Enhanced Nathan Kline Institute-Rockland (NKI) project (http://fcon_1000.projects.nitrc.org/indi/enhanced/). This pilot study comprises multimodal imaging data and subject-specific covariates for 24 individuals. In particular, for each subject, the following data are available:

- Structural networks measuring, from diffusion tensor imaging, white matter fiber interconnections among brain regions of interest.
- Functional activity data measuring the dynamic activity of each brain region through changes in the blood oxygen level-dependent signal during resting state functional magnetic resonance imaging.
- Functional networks denoting regions' synchronization in brain activity.

Spatial information on the brain regions of interest and subject-specific data on age, handedness, and psychological traits are also provided. The imaging data were pre-processed and generously provided by Greg Kiar and Eric Bridgford from NeuroData at Johns Hopkins University, who are gratefully acknowledged.

Motivated by the above dataset, the groups proposed stimulating methods during *StartUp Research* and continued their studies in the following year. More specifically, the contribution “[Understanding Dependency Patterns in Structural and Functional Brain Connectivity Through fMRI and DTI Data](#)” leverages latent variable models and dynamic Bayesian networks to learn, possibly similar, patterns in brain structural and functional connectivity. The contribution “[Hierarchical Graphical Model for Learning Functional Network Determinants](#)” instead adopts a modular approach which combines smoothing procedures, graphical models, and regression methods to relate functional connectivity with regions and subject-specific features. Different directions in the analysis of brain interconnections are proposed in “[Three Testing Perspectives on Connectome Data](#)”. The first focuses on learning structural restrictions in brain functional activity. The second aims at estimating the effective number of white matter fibers via parsimonious models, while the third studies group differences in brain connectivity with subjects’ traits under an object-oriented perspective. Also the work “[An Object Oriented Approach to Multimodal Imaging Data in Neuroscience](#)” analyzes the human brain data as object-valued and provides a wide set of procedures, including clustering, low-dimensional embeddings, and hypothesis testing, to obtain coherent findings in neuroscience. In a similar research direction, the contribution “[Curve Clustering for Brain Functional Activity and Synchronization](#)” focuses on appropriate methods to infer grouping structures and functional outliers in fMRI trajectories. These data are further explored in “[Robust Methods for Detecting Spontaneous Activations in fMRI Data](#)” via novel filtering methods which incorporate heavier tails than classical Gaussian assumptions. Parsimonious, yet flexible, Bayesian dynamic latent factor models are instead considered in the contribution “[Hierarchical Spatio-Temporal Modeling of Resting State fMRI Data](#)” to infer spatial and temporal effects of brain functional activity across multiple regions. A final article by Michele Guindani and Marina Vannucci summarizes the different proposals and opens toward new stimulating research directions in this field.

We would like to thank the early-career participants, Emanuele Aliverti, Gaia Bertarelli, Alessandra Cabassi, Alessia Caponera, Andrea Capozzo, Alessandro Casa, Alice Corbella, Federico Crescenzi, Marta Crispino, Silvia D’Angelo, Francesco Denti, Jacopo Di Iorio, Roberta Falcone, Federico Ferraccioli, Matteo Fontana, Laura Forastiere, Francesca Gasperoni, Anastasiia Gorshechnikova, Tullia Padellini, Sally Paganin, Michele Peruzzi, Alexios Polymeropoulos, Saverio Ranciati, Tommaso Rigon, Dutta Ritabrata, Massimiliano Russo, Andrea Sottosanti, and Marco Stefanucci for their enthusiasm and dedication to this stimulating experience. Also, the group leaders Alessio Farcomeni, Alan Gelfand, Alessandra Luati, Antonietta Mira, Piercesare Secchi, Marian Scott, and Ernst Wit are warmly acknowledged for their fundamental and inspiring contribution in leading the groups, both personally and scientifically.

Finally, we would like to thank the Italian Statistical Society, the Department of Economics and Statistics of the University of Siena, and the Department of Statistical Sciences of the University of Bologna for supporting *StartUp Research*. We are also grateful to the referees for their thoughtful revisions.

Padua, Italy

Milan, Italy

Milan, Italy

Padua, Italy

July 2018

Antonio Canale

Daniele Durante

Lucia Paci

Bruno Scarpa

Contents

Understanding Dependency Patterns in Structural and Functional Brain Connectivity Through fMRI and DTI Data	1
Marta Crispino, Silvia D'Angelo, Saverio Ranciati and Antonietta Mira	
Hierarchical Graphical Model for Learning Functional Network Determinants	23
Emanuele Aliverti, Laura Forastiere, Tullia Padellini, Sally Paganin and Ernst Wit	
Three Testing Perspectives on Connectome Data	37
Alessandra Cabassi, Alessandro Casa, Matteo Fontana, Massimiliano Russo and Alessio Farcomeni	
An Object Oriented Approach to Multimodal Imaging Data in Neuroscience	57
Andrea Cappozzo, Federico Ferraccioli, Marco Stefanucci and Piercesare Secchi	
Curve Clustering for Brain Functional Activity and Synchronization	75
Gaia Bertarelli, Alice Corbella, Jacopo Di Iorio, Anastasia Gorshechnikova and Marian Scott	
Robust Methods for Detecting Spontaneous Activations in fMRI Data	91
Francesca Gasperoni and Alessandra Luati	
Hierarchical Spatio-Temporal Modeling of Resting State fMRI Data	111
Alessia Caponera, Francesco Denti, Tommaso Rigon, Andrea Sottosanti and Alan Gelfand	
Challenges in the Analysis of Neuroscience Data	131
Michele Guindani and Marina Vannucci	

About the Editors

Antonio Canale is an Assistant Professor of Statistics in the Department of Statistical Sciences, University of Padova (Italy). His research areas cover Bayesian non-parametric methods, functional data analysis, statistical learning, and data mining. He is the author of a number of papers on methodological and applied statistics and has served on the scientific committees of national and international conferences. He was the coordinator of the young group of the Italian Statistical Society (y-SIS) in 2015.

Daniele Durante is an Assistant Professor of Statistics in the Department of Decision Sciences, Bocconi University (Italy), and a Research Affiliate at the Bocconi Institute for Data Science. His research is characterized by an interdisciplinary approach at the intersection of Bayesian methods, modern applications, and statistical learning to develop flexible and computationally tractable models for complex data. He is the coordinator of the young group of the Italian Statistical Society (y-SIS).

Lucia Paci is an Assistant Professor of Statistics in the Department of Statistical Sciences, Università Cattolica del Sacro Cuore, Milan (Italy). Her research focuses mainly on spatial and spatiotemporal modeling under the Bayesian framework, with applications in the environmental and economic sciences. She was the coordinator of the young group of the Italian Statistical Society (y-SIS) in 2016.

Bruno Scarpa is an Associate Professor of Statistics in the Department of Statistical Sciences, University of Padova (Italy). He teaches data mining at the master level and statistical methods for big data at the undergraduate level. His research interests include methodological developments motivated by real data applications. He is the author and coauthor of numerous papers and books in the fields of methodological and applied statistics and data mining.

Understanding Dependency Patterns in Structural and Functional Brain Connectivity Through fMRI and DTI Data



Marta Crispino, Silvia D'Angelo, Saverio Ranciati and Antonietta Mira

Abstract Neuroscience and neuroimaging have been providing new challenges for statisticians and quantitative researchers in general. As datasets of increasing complexity and dimension become available, the need for statistical techniques to analyze brain related phenomena becomes prominent. In this paper, we delve into data coming from functional Magnetic Resonance Imaging (fMRI) and Diffusion Tensor Imaging (DTI). The aim is to combine information from both sources in order to learn possible patterns of dependencies among regions of interest (ROIs) of the brain. First, we infer positions of these regions in a latent space, using the observed structural connectivity provided by the DTI data, to understand if physical spatial coordinates suitably reflect how ROIs are effectively interconnected. Secondly, we inspect Granger causality in the fMRI data in order to capture patterns of activations between ROIs. Then, we compare results from the analysis on these datasets, to find a link between functional and structural connectivity. Preliminary findings show that latent space positions well reflect hemisphere separation of the brain but are not perfectly connected to all the other structural partitions (that is, lobe, cortex, etc.); furthermore, activations of ROIs inferred from fMRI data are tied to observed structural connections derived from DTI scans.

M. Crispino
Univ. Grenoble Alpes, Inria, CNRS, LJK, 38000 Grenoble, France
e-mail: marta.crispino@inria.fr

S. D'Angelo
Department of Statistical Sciences, Sapienza University of Rome, Rome, Italy
e-mail: silvia.dangelo@uniroma1.it

S. Ranciati (✉)
Department of Statistical Sciences, University of Bologna, Bologna, Italy
e-mail: saverio.ranciati2@unibo.it

A. Mira
Institute of Computational Science, Università della Svizzera italiana, Lugano, Switzerland
e-mail: antonietta.mira@usi.ch

A. Mira
Department of Science and High Technology, Università dell'Insubria, Como, Italy

© Springer Nature Switzerland AG 2018

A. Canale et al. (eds.), *Studies in Neural Data Science*, Springer Proceedings in Mathematics & Statistics 257, https://doi.org/10.1007/978-3-030-00039-4_1

Keywords Network analysis · Resting state fMRI · DTI · Latent space models
Penalized weighted regression

1 Motivating Real World Dataset

Advances in neuroimaging have led to an increase in the availability of data to study complex systems (for instance, neurological processes in human brain). It is now possible to collect data considering different aims and assimilating different sources, a strategy that better captures the underlying dynamics of the phenomenon at study. The main interest lies in unraveling the mechanisms originating structural and functional brain activity, and simultaneously in understanding how these aspects are intertwined with patients covariates: for instance, how significant and relevant are the differences in brain connectivity and activity among subjects with heterogeneous characteristics. From different available multimodal brain imaging frameworks, here we focus on functional Magnetic Resonance Imaging (fMRI), and Diffusion Tensor Imaging (DTI).

The datasets we consider for the analysis were collected during a pilot study of the Enhanced Nathan Kline Institute-Rockland Sample project (information about the project itself is available at http://fcon_1000.projects.nitrc.org/indi/enhanced/). Data consist of 24 subjects, whose brain activity and structural connectivity were captured through DTI and resting state fMRI scan. The raw data were preprocessed and the scanned areas of the brain were parceled to determine a set of regions of interest (ROIs). An overview of the preprocessing steps is given in [4, 29]. Interest in analyzing rich and significant data from neuroimaging received a huge boost in the last decades, through a significant spillover of network analysis into neuroscience [2]. In an attempt to both distill information from complex systems and to infer the main mechanisms underlying brain activity, methods and concepts from network analysis were used into the framework of brain data. Concepts and terms such as hubs, centrality, hierarchy, node connectedness, and so forth, became both vocabulary and methodological tools shared by network analysis and neuroscience communities. A comprehensive review of this bridging between network analysis and neuroimaging is provided in [2]; more recently, in [30] an overview on connection between network properties and brain imaging data is also discussed, with emphasis on detection of neurological disorders via changes in the network structure itself. In [25, 29], Statistics and network science are directly tied to the study of functional and structural connectivity, and how both could help in deepening our understanding of interactions and, possibly, causality, among regions of the brain. In particular, for causal inference with an emphasis on fMRI data, refer to [22, 26].

From a statistical point of view, we find interest in exploiting—in a synergized approach—all the information at disposal. This can be done in many ways: for example, one could inject the number of white matter fibers for each pair of regions (DTI data) as a covariate information in a model capturing the correlation in the fMRI data. Alternatively, one might describe the statistical properties of an assumed underlying

network model originating the DTI dataset, then use some network's properties as aid in discovering (in fMRI data) patterns of synchronization between regions of the brain. Given the variety of questions and ways to address them, we adopt a statistical framework unifying these two ways of looking at brain connectivity, with the intent to assimilate different kinds of information captured by these different technologies, DTI and fMRI. The main goal of this work is thus to combine results from structural and functional observed data, in order to enhance the interpretation of each separate findings. In particular, our aim is to assess if the two datasets give 'coherent' answers with respect to the behavior we expect from the phenomenon: the patterns of activation among ROIs from fMRI data should be tied to structural connectedness highlighted by DTI data. The two datasets obviously share information on the overall activity of the brain but, from a modeling perspective, they bring different contributions to the whole research framework. In particular, data for the structural connectivity should depict the 'hardware' reference for us to understand which regions of the brain are physically connected. On the other hand, data from resting fMRI should provide insights on the dynamic counterpart of signals commuting between ROIs, and thus a different aspect on the concept of connectivity. For these reasons, we consider separate statistical models for the two available datasets.

First, we analyze the structural connectivity information provided by the DTI dataset, with the aid of models and tools coming from network analysis' framework. In particular, we investigate the idea that a statistical interpretation of the topology of the network differs from the physical observed topology, represented by the spatial coordinates of the ROIs. For this reason, we resort on latent space models which allow us to infer positions of the ROIs, in terms of how close they are, directly from the data on white matter fibers and their structural connectivity. Latent space models for social network analysis have been introduced in [15]. In their work, the authors assume that the observed network data depend on a set of latent variables. Indeed, the nodes are assumed to be in a p -dimensional latent space. Then, the probability that two nodes are joined by an edge in the network depends on some function of the unknown latent coordinates of the nodes. In the case of *distance latent space models*, this function is generally assumed to be the Euclidean distance: the smaller the distance, the greater the probability of an edge. While for *projection latent space model*, the function considers the angle formed between two nodes in the bilinear latent space: the smaller the angle the higher the probability that the dyad is connected. This class of models can take into account some of the typical features of network analysis, such as the presence of degree heterogeneity and of group structure. Indeed [14] introduced sender and receiver effects for networks and [13] proposed a clustering model for the nodes in the latent space. A different approach that makes use of latent variable to model the dependency structure observed in network data is *stochastic block modeling* [21, 27]. Stochastic block models are particularly suited to cluster the nodes into blocks. From this rich literature, we mainly draw from the contribution of [15] in order to gain insights about spatial organizations of the ROIs.

Second, we deal with temporal information provided by the dataset on resting state fMRI, which comprises of time series of brain activity. In particular, we exploit a weighted linear regression model that encodes a type of causality between

observations at different time points and for different ROIs. This causality, called *Granger causality*, is tied to the estimation of parameters in the weighted linear regression. Granger Causality [9, 10] was first introduced in the Econometric literature, specifically developed for time series analysis. This notion of causality is grounded on the rather obvious intuition that the origin of a cause should necessarily precede in time its effect. In particular, it states that, given two sets of time series data, V_1 and V_2 , the series V_1 (Granger) causes V_2 if past values of V_1 are helpful in predicting the future values of V_2 . It is important to underline that Granger causality is not intended as causality in a deep sense: it just measures whether one time series is likely to influence the other one, that is, if V_1 provides more information about future values of V_2 than past values of V_2 alone. As such, Granger causality not always overlaps with actual causality, but it is still a useful instrument to infer whether two series are related by some, generally unknown, phenomenon. Recently, the notion of Granger causality entered into the network literature on multivariate time series, with the objective of learning sparse sets of Granger causal relationships between univariate series [11, 12, 17, 24, 31]. In this paper, we exploit this notion in order to infer any existing relationship between different brain regions.

In light of the above discussion, the remainder of the paper is organized as follows: in Sect. 2 we provide some exploratory statistics on the datasets, to summarize the salient features of the data we are modeling; in Sect. 3, we outline a modeling approach to static network data along with the results we get by applying this methodology to the DTI dataset (Sect. 3.1). Then, a time-varying dynamic linear model formulation for the time-series dataset is presented in Sect. 4, together with preliminary results from the fMRI dataset (Sect. 4.1). Finally, in Sect. 5, we discuss the results so far obtained, also providing a glimpse of future developments.

2 Descriptive Analysis

As mentioned in Sect. 1, we here analyze two core datasets. The first one, DTI, refers to subjects' brain structural connectivity, measured using DTI. For the scan of each individual (and re-scan, if available), a 70×70 matrix reports the observed count of white matter fibers connecting pairs of ROIs; a structural 'NA' (not assigned) value is reported for self-connectivity, which is the diagonal of the aforementioned matrix. The number of white matter fibers is thus an observed measure of connectivity between brain regions. To compare our results with a standard brain representation, we refer to the atlas for brain parcellation in [6] (also reported in Appendix 6), that has a total of $n = 68$ ROIs. To match the analyzed regions with those of the Desikan atlas, two of the 70 ROIs in the data labeled as "unknown" are discarded, when comparing the results. Moreover, two of the regions in the data refer to the corpus callosum (left and right), while the Desikan atlas does not represent this region but instead considers the insula, which includes the corpus callosum together with the lateral ventricles. Therefore, when comparing our results with the Desikan atlas representation of the brain, we refer to the corpus callosum as insula. The

second dataset, fMRI, comprises of resting state dynamic functional activities of the ROIs, measured via blood-oxygen-level dependent (BOLD) technique through fMRI, at $T = 404$ equally spaced time points with in-between lags of 1400 (ms). A third collateral dataset was produced by computing, from the fMRI dataset, time-wise correlations among the ROIs, resulting in a 70×70 matrix for each subject (and each re-scan, when available). Additional information is provided in the form of covariates. For the subjects, characteristics available are: status of current (single episode/recurrent depressive disorder, cannabis abuse, anxiety, social phobia) or lifetime (alcohol abuse, drug dependence, Attention Deficit Hyperactivity Disorder, eating disorder, major depressive disorder) mental disorder; handedness (left, right, ambidextrous); age. For the ROIs, their lobe and hemisphere memberships are recorded, together with the physical spatial coordinates of the centroids used in the atlas. Information is not always available for all individuals and, in some cases, not every subject has a re-scan dataset to be paired with the original scan, causing missing values in the reported observations.

We here provide some descriptive statistics to familiarize the reader with the data at hand. In Table 1 we focus on the DTI dataset. The data for 4 subjects (labeled 6, 17, 20 and 22) are not available, and the disease diagnosis is missing for four patients (labeled 3, 4, 5, 6). We also notice that the seven patients who are diagnosed with lifetime disease, are also diagnosed with current disease. More importantly, three out of four patients with a diagnosis of current disease have missing data. The descriptive analysis in this section therefore focuses only on differences among patients with diagnosis of lifetime disease (four available out of seven), and patients with no such diagnosis (twelve available out of thirteen). The most salient feature of Table 1 is that the variability within subject of the number of white matter fibers is extremely high, with a range of values between zero and several thousands; distribution of the number of white matter fibers for each patient is highly skewed, and the median values are always much smaller than the mean values. We check if adjusting marginally (that is, one at a time) for covariates can help explain the distribution of the median number of white matter fibers across the patients. For example, the boxplot in Fig. 1 (left panel) represents the median number of white matter fibers, stratified by the lifetime disease diagnosis of the subjects (YES/NO).

From this plot, we notice that there is a difference, in terms of median number of white matter fibers, among patients with and without lifetime disease diagnosis. This result is confirmed by a Wilcoxon rank sum test which rejects the null hypothesis (p -value = 0.0077) of equal medians in the two groups. We also performed a Wilcoxon rank sum test to assess whether there is a difference between subjects with/without diagnosis of lifetime disease in terms of the percentage of zeros in the adjacency matrix (column 8 of Table 1). The null hypothesis in this case is not rejected. These findings suggest that the number of white matter fibers may be particularly informative about disease status, while the presence/absence of white matter fibers is not. There is a huge literature (see e.g [18, 20, 23]) that studies associations between mental disorders and DTI of white matter fibers. Analogous considerations can be

Table 1 Summary table of the DTI dataset. Grey rows refer to patients who received a lifetime diagnosis. Empty cells correspond to non-available data. Columns (Statistics) from left to right: minimum value, maximum value, mean, median, standard deviation, interquartile range, coefficient of variation, percentage of zeros in the adjacency matrix

	Statistics					Covariates						
	Min	Max	Mean	Median	sd	iqr	cv	% of 0s	Age	Handed	Current diagnosis	Lifetime diagnosis
	(computed on the non-0 entries)											
Patient 1	1	29078	2340	789	3814	2709	163	0.53	57	R	NO	NO
Patient 2	1	38041	2802	1009	4562	3315	163	0.63	52	R		
Patient 3	1	33081	2978	968	4822	3389	162	0.62	32	R		
Patient 4	1	39783	3068	1093	4771	4218	155	0.57	36	R		
Patient 5	1	34807	2840	995	4426	3428	156	0.67	22	R		
Patient 6									27	R	NO	NO
Patient 7	1	34481	2844	927	4738	3010	167	0.65	60	R	NO	YES
Patient 8	1	36809	3375	1057	5351	4101	159	0.61	21	R	YES	YES
Patient 9	1	37286	2566	898	4207	3019	164	0.60	21	L	NO	NO
Patient 10	1	33858	3123	1058	4949	3690	158	0.62	30	R	NO	NO
Patient 11	1	31682	2198	492	3848	2634	175	0.53	27	A	NO	NO
Patient 12	1	29288	2832	1005	4383	3260	155	0.63	48	A	NO	NO
Patient 13	1	41524	3530	1406	5340	4508	151	0.60	22	R	NO	YES
Patient 14	1	32377	2805	973	4446	3461	158	0.63	19	R	NO	NO
Patient 15	1	26388	2499	792	3826	3238	153	0.64	57	R	NO	NO
Patient 16	1	28797	2251	672	3484	3043	155	0.65	25	R	NO	NO
Patient 17									38	R	YES	YES
Patient 18	1	27501	1988	692	3070	2556	154	0.66	46	R	NO	NO
Patient 19	1	44968	3645	1179	6037	4500	166	0.62	22	R	NO	YES
Patient 20									32	L	YES	YES
Patient 21	1	36270	2888	791	4839	3627	168	0.60	22	R	NO	NO
Patient 22									42	R	YES	YES
Patient 23	1	36857	2585	813	4191	3095	162	0.60	31	L	NO	NO
Patient 24	1	41699	2990	1050	4862	3441	163	0.63	36	R	NO	NO

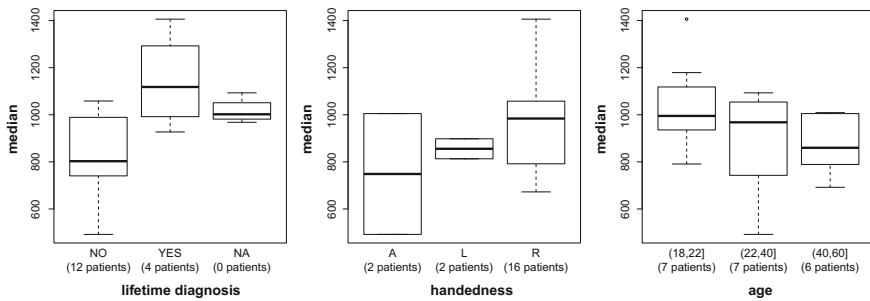


Fig. 1 Boxplots of the median number of white matter fibers stratified by the lifetime disease diagnosis (left), handedness (middle), age (right)

drawn for other covariates, such as age or handedness¹ (see Fig. 1, middle and right panels). However, in this paper we do not model directly the number of white matter fibers (see Sect. 3), or the potential impact of covariates on these counts. Rather, we look at absence/presence of fibers since our focus is mainly in joining information coming from the two datasets. Nevertheless, we are willing to investigate this aspect in a future development of the current analysis (see also Sect. 5).

The *fMRI* dataset reports the dynamic activity time series of each brain region, for each subject in the study. Data for 2 subjects (1 and 21) are not available, whereas for other 11 subjects (labeled 4, 10, 11, 12, 13, 14, 15, 16, 17, 18, 20) only the first scan is observed. In Fig. 2 we report some plots produced as follows:

- for each subject and each region of the brain, we compute the range of the activation levels along the time series;
- we then plot the obtained values versus the regions on the *x*-axis, with different symbols related to the user-specific covariates (as in the legend).

The range for patients who are diagnosed with a lifetime disease (black squares in the upper panel of Fig. 2) seems to be much higher than the one of the patients who did not receive a positive diagnosis (grey circles). Differences are noticeable also inspecting the handedness plot, in the lower panel of Fig. 2. We see that the range of activation is smaller for ambidextrous (A) patients (light-grey triangles).

The third dataset reports synchronization in brain activity for each pair of brain regions, obtained from the correlation in the dynamic functional activity of the *fMRI* dataset. In the heatmap of Fig. 3 is depicted the adjacency matrix corresponding to the functional network measuring correlation in brain activity between pairs of regions. The upper triangular panel is built averaging values of the patients who were not diagnosed a lifetime disease, while the lower triangular one refers to patients with lifetime disease. The structure of the two triangular panels is very similar, as indicated by the dark cells common in the two panels, and by the evident subdivision of both plots into three areas. However, it is evident an overall higher correlation for patients with diagnosis, in particular regarding pairs of brain regions that do not activate

¹Handedness is the dominance of one hand over the other, or the unequal distribution of fine motor skills between the left and right hands.