



**SDS**  
Statistica e  
Data Science



Palermo, 11-12 April 2024

# Proceedings of the Statistics and Data Science 2024 Conference

---

## New perspectives on Statistics and Data Science

Edited by

---

Antonella Plaia – Leonardo Egidi  
Antonino Abbruzzo

Proceedings of the SDS 2024 Conference  
Palermo, 11-12 April 2024  
Edited by: Antonella Plaia - Leonardo Egidi - Antonino  
Abbruzzo

-

Palermo: Università degli Studi di Palermo.

ISBN Ebook (Pdf)  
978-88-5509-645-4

Questo volume è rilasciato sotto licenza Creative Commons  
**Attribuzione - Non commerciale - Non opere derivate 4.0**



© 2024 The Authors

# Contents

<b>1</b>	<b>Keynote Sessions</b> .....	<b>11</b>
1.1	Classification with imbalanced data and the (eternal?) struggle between statistics and data science. <i>Nicola Torelli</i>	11
1.2	Deep residual networks and differential equations. <i>G�rard Biau</i>	13
<b>2</b>	<b>Invited - Complex data: new methodologies and applications</b> ..	<b>15</b>
2.1	Link selection in binary regression models with the Model Confidence Set. <i>Michele La Rocca and Marcella Niglio and Marialuisa Restaino</i>	15
2.2	A cluster-weighted model for COVID-19 hospital admissions. <i>Daniele Spinelli, Paolo Berta, Salvatore Ingrassia and Giorgio Vittadini</i>	23
2.3	Multi-class text classification of news data. <i>Maurizio Romano and Maria Paola Priola</i>	28
<b>3</b>	<b>Invited - Data science and dataspace: challenges, results, and next steps</b> .....	<b>35</b>
3.1	Data-Centric AI : A new Frontier emerging in Data Science. <i>Donato Malerba, Vincenzo Pasquadibisceglie, Vito Recchia and Annalisa Appice</i>	35
3.2	Data Spaces strategy to unleash agriculture data value: a concrete use case. <i>Nicola Masi, Delia Milazzo, Giulia Antonucci and Susanna Bonura</i>	42
3.3	Addressing Agricultural Data Management Challenges with the Enhanced TRUE Connector. <i>Sergio Comella, Delia Milazzo, Mattia Giuseppe Marzano, Giulia Antonucci, Susanna Bonura and Angelo Marguglio</i>	48
<b>4</b>	<b>Solicited - Data Science for Official Statistics</b> .....	<b>55</b>
4.1	Data science at Istat for urban green. <i>Fabrizio De Fausti, Marco Di Zio, Giuseppe Lancioni, Stefano Mugnoli, Alberto Sabbi and Francesco Sisti</i>	55

4.2	Twitter (X) as a Data Source for Official Statistics: Monitoring Italian Debate on Immigration through Text Analysis. <i>Elena Catanese, Gerarda Grippo, Francesco Ortame and Maria Clelia Romano</i>	62
<b>5</b>	<b>Solicited - Sustainable Artificial Intelligence in Finance</b> .....	<b>69</b>
5.1	Feature Dependence and Prediction Explanations in P2P Lending. <i>Paolo Pagnottoni and Thanh Thuy Do</i>	69
<b>6</b>	<b>Solicited - Young SIS</b> .....	<b>77</b>
6.1	Merging data and historical information via optimal power prior selection in Bayesian models. <i>Roberto Macrì Demartino, Leonardo Egidi, Nicola Torelli and Ioannis Ntzoufras</i>	77
6.2	Hierarchical Mixtures of Latent Trait Analyzers with concomitant variables. <i>Dalila Failli, Bruno Arpino, and Maria Francesca Marino</i>	84
6.3	A Simultaneous Spectral Clustering for Three-Way Data. <i>Cinzia Di Nuzzo and Salvatore Ingrassia</i>	90
<b>7</b>	<b>Solicited - From Data Analysis to Data Science</b> .....	<b>97</b>
7.1	Optimal Scaling: New Insights Into an Old Problem. <i>Gilbert Saporta</i>	97
<b>8</b>	<b>Solicited - Statistical methods for textual data</b> .....	<b>101</b>
8.1	PROCSIMA: Probability Distribution Clustering Using Similarity Matrix Analysis. <i>Marco Ortu</i>	101
8.2	Exploring Anti-Migrant Rhetoric on Italian Social Media. <i>Lara Fontanella, Annalina Sarra, Emiliano del Gobbo, Alex Cucco and Sara Fontanella</i>	108
8.3	Causal inference from texts: a random-forest approach. <i>Chiara Di Maria, Alessandro Albano, Mariangela Sciandra and Antonella Plaia</i>	114
<b>9</b>	<b>Solicited - Data analysis methods for data in non-Euclidean spaces</b>	<b>121</b>
9.1	Riemannian Statistics for Any Type of Data. <i>Oldemar Rodriguez Rojas</i>	121
9.2	PAM clustering algorithm for ATR-FTIR spectral data selection: an application to multiple sclerosis. <i>Francesca Condino, Maria Caterina Crocco and Rita Guzzi</i>	128
9.3	Random Survival Forest for Censored Functional Data. <i>Giuseppe Loffredo, Elvira Romano and Fabrizio Maturo</i>	134
9.4	Advancing credit card fraud detection with innovative class partitioning and feature selection technique. <i>Mohammed Sabri, Antonio Balzanella and Rosanna Verde</i>	140
<b>10</b>	<b>Solicited - Functional Data Analysis in Action</b> .....	<b>147</b>
10.1	Functional Linear Discriminant Analysis for Misaligned Surfaces. <i>Tomas Masak</i>	147
10.2	Leveraging weighted functional data analysis to estimate earthquake-induced ground motion. <i>Teresa Bortolotti, Riccardo Peli, Giovanni Lanzano, Sara Sgobba and Alessandra Menafoglio</i>	155

10.3	Functional autoregressive processes on a spherical domain for global aircraft-based atmospheric measurements. <i>Almond Stöcker and Alessia Caponera</i>	161
<b>11</b>	<b>Solicited - Bayesian Inference for Graphical Models</b> .....	<b>169</b>
11.1	Log-likelihood approximation in Stochastic EM for Multilevel Latent Class Models. <i>Silvia Columbu, Nicola Piras and Jeroen K. Vermunt</i>	169
11.2	MCMC Sampling in Bayesian Gaussian Structure Learning. <i>Antonino Abbruzzo, Nicola Argentino, Reza Mohammadi, Maria De Iorio, Willem van den Boom and Alexandros Beskos</i>	176
<b>12</b>	<b>Contributed - Promoting Equity: Statistical Insights into Tourism, Sustainability and Digital Divide</b> .....	<b>183</b>
12.1	Lesson Learnt in the Data Science Worldview: New Dimension of Digital Divide. <i>Rita Lima</i>	183
12.2	An overview of Tourism Statistical Literacy. <i>Yasir Jehan, Giuseppina Lo Mascolo and Stefano De Cantis</i>	192
12.3	Scalable bootstrap inference via averaged Robbins-Monro approximations. <i>Giuseppe Alfonzetti and Ruggero Bellio</i>	198
12.4	The impact of sustainability on Initial Coin Offering: advantages in trading. <i>Alessandro Bitetto and Paola Cerchiello</i>	204
<b>13</b>	<b>Contributed - High dimensional and functional data</b> .....	<b>211</b>
13.1	Analysis of Brain-Body Physiological Rhythm Using Functional Graphical Models. <i>Rita Fici, Luigi Augugliaro and Ernst C. Wit</i>	211
13.2	A comparison of scalable estimation methods for large-scale logistic regression models with crossed random effects. <i>Ruggero Bellio and Cristiano Varin</i>	218
13.3	Single-cell Sequencing Data: Critical Analysis and Definition of Statistical Models. <i>Antonino Gagliano, Gianluca Sottile, Nicolina Sciaraffa, Claudia Coronello and Luigi Augugliaro</i>	224
13.4	Investigating the association between high school outcomes and university enrollment choice: a machine learning approach. <i>Andrea Priulla, Alessandro Albano, Nicoletta D'Angelo and Massimo Attanasio</i>	230
<b>14</b>	<b>Contributed - Statistical Analysis in economic and market dynamics</b>	<b>237</b>
14.1	A comparison of multi-factor stochastic models for commodity prices C3. <i>Luca Vincenzo Ballestra, Christian Tezza and Paolo Foschi</i>	237
14.2	Nonparametric ranking estimation with application to the propensity for Circular Economy of Italian economic sectors. <i>Stefano Bonnini, Michela Borghesi and Massimiliano Giacalone</i>	246
14.3	Impact of the Russian invasion of Ukraine on coal markets: Evidence from an event-study approach. <i>Yana Kostiuk, Paola Cerchiello and Arianna Agosto</i>	252
14.4	Labour market and time series: a forecast approach for European countries from 1995 to 2022. <i>Paolo Mariani, Andrea Marletta and Piero Quatto</i>	258

<b>15</b>	<b>Contributed - Innovations in cluster and latent class models</b> . . .	<b>263</b>
15.1	Biclustering of discrete data by extended finite mixtures of latent trait models. <i>Dalila Failli, Maria Francesca Marino and Francesca Martella</i>	263
15.2	Seismic events classification through latent class regression models for point processes. <i>Giada Lo Galbo, Giada Adelfio and Marcello Chiodi</i>	270
15.3	Determining the optimal number of clusters through Symmetric Non-Negative Matrix Factorization. <i>Agostino Stavolo, Maria Gabriella Grassia, Marina Marino and Rocco Mazza</i>	276
<b>16</b>	<b>Contributed - Modelling on spatial phenomena</b> . . . . .	<b>283</b>
16.1	Integrating computational and statistical algorithms in RT-GSCS for spatial survey administration. <i>Yuri Calleo, Simone Di Zio and Francesco Pilla</i>	283
16.2	Sensitivity mapping as a tool to support siting of offshore wind farms and increase citizens' acceptability. <i>Giovanna Cilluffo, Gianluca Sottile, Laura Ciriminna, Geraldina Signa, Agostino Tomasello and Salvatrice Vizzini</i>	290
16.3	Investigating hotel consumers' purchase intention on web analytics data through PLS-SEM. <i>Giuseppina Lo Mascolo, Chiara di Maria, Marcello Chiodi and Arabella Mocciano Li Destri</i>	296
16.4	Spatio-temporal analysis of lightning point process data in severe storms. <i>Nicoletta D'Angelo, Milind Sharma, Marco Tarantino and Giada Adelfio</i>	302
<b>17</b>	<b>Contributed - Statistical machine learning for predictive modelling</b>	<b>309</b>
17.1	Application of statistical techniques to predict the effective temperature of young stars. <i>Marco Tarantino, Loredana Prisinzano and Giada Adelfio</i>	309
17.2	Topological Attention for Denoising Astronomical Images. <i>Riccardo Cecaroni and Pierpaolo Brutti</i>	316
17.3	LSTM-based Battery Life Prediction in IoT Systems: a proof of concept. <i>Vanessa Verrina, Andrea Vennera and Annarita Renda</i>	322
17.4	Predictive modeling of drivers' brake reaction time through machine learning methods. <i>Alessandro Albano, Giuseppe Salvo and Salvatore Rusotto</i>	328
<b>18</b>	<b>Contributed - Ordinal and preference data analysis</b> .	<b>335</b>
18.1	OSILA (Order Statistics In Large Arrays): an original algorithm for an efficient attainment of the order statistics. <i>Andrea Cerasa</i>	335
18.2	The Mallows model with respondents' covariates for the analysis of preference rankings. <i>Marta Crispino, Lucia Modugno and Cristina Mollica</i>	343
18.3	Value-Based Predictors of Voting Intentions: An Empirical and Explainable approach. <i>Luca Pennella and Amin Gino Fabbrucci Barbagli</i>	349
18.4	A dynamic version of the Massey's rating system with an application in basketball. <i>Paolo Vidoni and Enrico Bozzo</i>	355
<b>19</b>	<b>Contributed - Advances in text mining</b> . . . . .	<b>361</b>
19.1	Can Correspondence Analysis Challenge Transformers in Authorship Attribution Tasks?. <i>Andrea Sciandra and Arjuna Tuzzi</i>	361

- 19.2 A Fuzzy Topic Modeling approach to legal corpora. *Antonio Calcagni and Arjuna Tuzzi* 368
- 19.3 EmurStat: a digital tool for statistical analysis of emur flow. *Simone Paesano, Maria Gabriella Grassia, Marina Marino, Dario Sacco and Rocco Mazza* 374
- 19.4 Graph Neural Networks for clustering medical documents. *Vittorio Torri and Francesca Ieva* 380





## Preface

The development of large-scale data analysis and statistical learning methods for data science is gaining more and more interest, not only among statisticians, but also among computer scientists, mathematicians, computational physicists, economists, and, in general, all experts in different fields of knowledge who are interested in extracting insight from data. Cross-fertilization between the different scientific communities is becoming crucial for progressing and developing new methods and tools in data science. In this respect, the Statistics & Data Science group of the Italian Statistical Society has organized its 3rd international conference held in Palermo on the 11st and 12nd of April 2024, attended by over 100 researchers from different scientific fields. A collection of the presented papers is available in the present Proceedings showing a huge variety of approaches, methods, and data-driven problems, always tackled according to a rigorous and robust scientific paradigm.

The Statistics & Data Science group

*Palermo, April 11st and 12th, 2023*

*Antonella Plaia - Leonardo Egidi - Antonino Abbruzzo*

*Editors*

## **13.** Contributed - High dimensional and functional data

13.1 - Analysis of Brain-Body Physiological Rhythm Using Functional Graphical Models

13.2 - A comparison of scalable estimation methods for large-scale logistic regression models with crossed random effects

13.3 - Single-cell Sequencing Data: Critical Analysis and Definition of Statistical Models

13.4 - Investigating the association between high school outcomes and university enrolment choice: a machine learning approach

# Analysis of Brain–Body Physiological Rhythm Using Functional Graphical Models

Fici R., Augugliaro L. and Wit E.C.

**Abstract** This paper presents an analysis of physiological data derived from a recent investigation on network physiology, adopting the conceptual framework that views the human organism as a complex network of interacting organs. The study explores coordinated interactions among organs using functional conditional Gaussian Graphical Models (fcGGM). Organ functions are modelled as networks with individual regulatory mechanisms, forming a broader system through continuous interactions. The focus of the analysis is on the interactions within and between two subnetworks: brain activity and a composite network comprising the RR interval of the electrocardiographic waveform, respiration amplitude and blood volume pulse.

**Key words:** EEG waves, network physiology, Functional Data Analysis, conditional functional Graphical Models

## 1 Introduction

The paper considers the intricate dynamics of physiological interactions, offering insights into the temporal synchronization of organ functions at the whole organism level, using functional data analysis (FDA). This analysis contributes to exploring new statistical tools to advance our comprehension of physiology interaction, shedding light on the intricate mechanisms that underlie the orchestration of physiological states within the human organism. Adopting a functional framework, the

---

Rita Fici  
University of Palermo, Viale delle Scienze - Building 13, Palermo, e-mail: rita.fici@unipa.it

Luigi Augugliaro  
University of Palermo, Viale delle Scienze - Building 13, Palermo, e-mail: luigi.augugliaro@unipa.it

Ernst-Jan Camiel Wit  
Università della Svizzera italiana, Via Buffi 13, Lugano, e-mail: wite@usi.ch

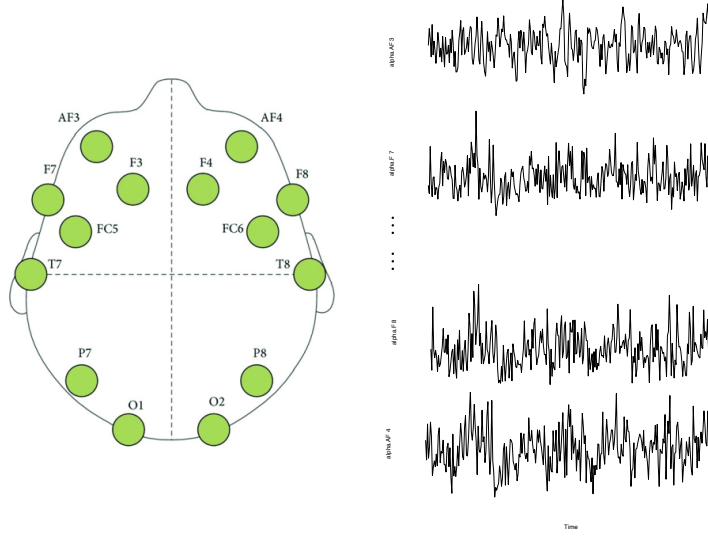
vector of records for each variable and each statistical unit represents a discretized observation of a functional variable. Both the response and the explanatory sets of variables constitute two functional processes which take place in an interval of the real line  $\mathcal{T} \in \mathbb{R}$  and denoted by  $\mathbf{Y}(t)$  and  $\mathbf{X}(t)$ , where  $t \in \mathcal{T}$ . This analysis aims to recover the conditional independent structure of brain waves and the effect of some external covariates on the expected function of each brain wave. These relations can be represented by a graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ , where  $\mathcal{V}$  is the set of variables given by the response and the explanatory functions, and  $\mathcal{E}$  represents undirected edges between functions. We apply the conditional functional Gaussian Graphical model proposed by [1], where the cross-covariance operator is considered as in [4] and the functional processes are assumed to have a partially separable covariance operator defined in [7].

## 2 Analysis

In this section, we first provide details on the analyzed data. Then, in Section 2.2, we offer a brief overview of the applied model and the estimation procedure used. Lastly, in Section 2.3, the estimation results are presented.

### 2.1 Brain-Body data

The data utilized in this study originate from a recent investigation on network physiology [3]. Within their conceptual framework, the human organism is viewed as a network where coordinated interactions among organs are essential for generating distinct physiological states and maintaining health. Each organ is considered as a network with its own regulatory mechanisms, continuously interacting with other body organs to contribute to the broader system. Physiological interactions occur at multiple levels across the temporal scale to synchronize organ functions at the whole organism level. Our analysis focuses on the interaction within and between two sub-networks;  $\mathbf{Y}(t)$  is the network consisting of the  $\alpha$ -waves of brain activity, and  $\mathbf{X}(t)$  describes the activity of three organs, namely: the time interval between two successive R-waves on the electrocardiographic waveform (RR), respiration amplitude, and blood volume pulse (BVP). Brain activity is recorded using the EPOC PLUS wireless headset, consisting of 14 electrodes that record signals from different brain regions. The ECG (sampled at 250 Hz) and respiration signals (sampled at 25 Hz) are acquired through a sensorized t-shirt, while the E4 wristband captures the BVP signal (sampled at 64 Hz) with a photoplethysmographic (PPG) sensor. Each electrode records the amplitude ( $\alpha$ ) of EEG waves. As described by [8], the data were collected by monitoring eighteen young, healthy volunteers during a measurement protocol consisting of three experimental conditions corresponding to different levels of mental stress. In this study, we analyze the data corresponding to a resting



**Fig. 1** Left: Graphical representation of the positioning of the 14 EEG electrodes over the scalp.

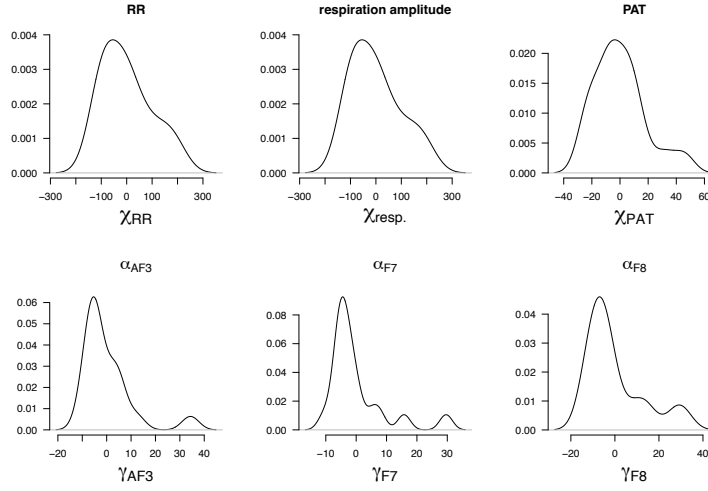
condition, lasting 12 min and consisting of watching a video showing landscapes with relaxing background music. The data consists of 300 time instants. Figure 1 shows the placement of the EEG electrodes over the scalp and some of the lines designed by the observed 300-dimensional discretized version of the brain waves.

## 2.2 Functional analysis methods

To recover the edge set representing the relation within and between the two processes, we analyze the relations among the coefficients resulting from the truncated joint vectorized Karhunen-Loève expansion (see [1]):

$$\mathbf{Z}(t) = \sum_{m=1}^M \boldsymbol{\zeta}_m \varphi_m(t)$$

where  $\mathbf{Z}(t) = [\mathbf{Y}(t)^\top \mathbf{X}(t)^\top]^\top$ , and each element in  $\boldsymbol{\zeta}_m$  is equal to  $\zeta_{m,i} = \int_{\mathcal{T}} Z_i(t) \varphi_m(t) dt$  where  $\{\varphi_m\}_{m=1}^M$  is a system of orthonormal basis. To estimate the basis, we perform the eigen-decomposition on the observed matrices of data (see [2]). We analyze  $\{\boldsymbol{\zeta}_m\}_{m=1}^{18}$ , which constitute 90% of the total variability. Figure 2 shows the kernel estimations of the densities of the scores for  $M = 1$  for some of the variables under examination. The assumption of Gaussianity underlying the fcMMs appears quite appropriate for the scores of the explanatory variables. However, an observed asymmetry is noted in the scores of the response variables. To recover the edge-set, we



**Fig. 2** Estimated densities for the scores of the explanatory variables and some of the response variables.

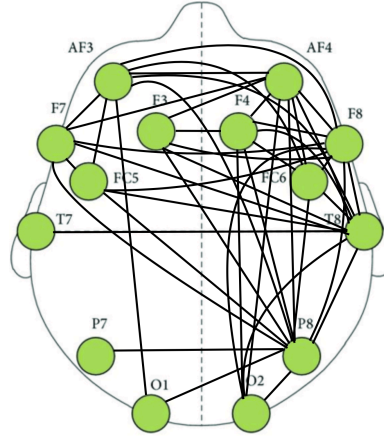
use the double-penalized joint group lasso estimator proposed by [5]. The optimization problem is performed iteratively, fixing one of the two tuning parameters and varying the other one with 20 values. The tuning parameter value for the next iteration is selected using the joint version of Kullback-Leibler cross-validation (KLCV, [6]) proposed in [1]. The total number of iterations is 3.

### 2.3 Results

Figure 3 shows the selected links representing the response conditional independence structure. The right side of the scalp exhibits a higher number of links compared to the left side. The most connected alpha wave is  $P8$  with 12 links. The waves  $T8$  and  $F8$  have 11 connection each. Table 1 provides details on the joint KLCV (jKLCV) regressors selection, indicating which explanatory variables have impacts on each of the response variables. The RR result affects all the  $\alpha$  waves. The respi-

	AF3	F7	F3	FC5	T7	P7	O1	O2	P8	T8	FC6	F4	F8	AF4
RR	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Resp	1	1	1	1	1	1	1	1	0	1	1	1	1	1
PAT	0	0	0	1	1	1	0	0	0	0	0	0	0	0

**Table 1** jKLCV: regressors selection as results of jKLCV selection method.



**Fig. 3** jKLCV: conditional independent response network.

ration amplitude affects all the response waves except *P8*. Lastly, the blood volume pulse affects the  $\alpha$  waves of *FC5*, *T7* and *P7* sensors.

### 3 Conclusion

The present work aims to contribute to the refinement and development of methods that allow addressing questions related to physiology by adopting a functional point of view.

This study delves into the intricate dynamics of physiological interactions, employing functional data analysis (FDA) to explore the temporal synchronization of organ functions at the whole organism level. By applying the conditional functional Gaussian Graphical model, we aimed to uncover the conditional independence structure of brain waves and the influence of external covariates. The adopted methodology includes the truncated joint vectorized Karhunen-Loève expansion and the double-penalized joint group lasso estimator. Those tools allow us to recover the conditional independence structure and identify influential factors on brain wave activities. The results highlighted the significance of the RR interval and respiration amplitude in shaping the expected value of almost all the  $\alpha$  brain waves. The blood volume pulse has a significant effect on some of the sensors.

## References

1. Fici, R., Augugliaro L., Wit E.C.: Functional Conditional Gaussian Graphical Models. (2024) doi: 10.48550/arXiv.2401.10196
2. Hsing, T., Eubank, R.: Theoretical foundations of functional data analysis, with an introduction to linear operators. John Wiley & Sons. (2015) doi: 10.1002/9781118762547
3. Pernice R, Antonacci Y, Zanetti M, Busacca A, Marinazzo D, Faes L, Nollo G.: Multivariate Correlation Measures Reveal Structure and Strength of Brain-Body Physiological Networks at Rest and During Mental Stress. *Front Neurosci. Frontiers in Neuroscience.* **14** (2021) doi: 10.3389/fnins.2020.602584
4. Qiao, X., Guo, S., James, G.M.: Functional Graphical Models. *Journal of the American Statistical Association.* **114** (525) pp. 211-222 (2019) doi: 10.1080/01621459.2017.1390466
5. Sottile G., Augugliaro L., Vinciotti V, Arancio W., Coronello C.: Sparse inference of the human hematopoietic system from heterogeneous and partially observed genomic data. (2022) doi: 10.48550/arXiv.2206.09863
6. Vujačić, I., Abbruzzo, A., Wit E.C.: A computationally fast alternative to cross-validation in penalized Gaussian graphical models. *Journal of Statistical Computation and Simulation.* **85** (18) pp. 3628-3640 (2015) doi: 10.1080/00949655.2014.992020
7. Zapata J., Oh S.Y., Petersen A.: Partial separability and functional graphical models for multivariate Gaussian processes. Oxford University Press (OUP). **109**(3) pp. 665–681 (2021) doi: 10.1093/biomet/asab046
8. Zanetti, M., Mizumoto, T., Faes L., Fornaser, A., Cecco, M., Maule, L. Valente, M., Nollo, G.: Multilevel assessment of mental stress via network physiology paradigm using consumer wearable devices. *Journal of Ambient Intelligence and Humanized Computing.* **12** pp. 4409–4418 (2021) doi: 10.3389/fnins.2020.602584