



# CLADAG 2021

BOOK OF ABSTRACTS AND SHORT PAPERS  
13th Scientific Meeting of the Classification and Data Analysis Group  
Firenze, September 9-11, 2021

edited by

Giovanni C. Porzio

Carla Rampichini

Chiara Bocci



PROCEEDINGS E REPORT

ISSN 2704-601X (PRINT) - ISSN 2704-5846 (ONLINE)

## SCIENTIFIC PROGRAM COMMITTEE

Giovanni C. Porzio (chair) (University of Cassino and Southern Lazio - Italy)

Silvia Bianconcini (University of Bologna - Italy)

Christophe Biernacki (University of Lille - France)

Paula Brito (University of Porto - Portugal)

Francesca Marta Lilja Di Lascio (Free University of Bozen-Bolzano - Italy)

Marco Di Marzio ("Gabriele d'Annunzio" University of Chieti-Pescara - Italy)

Alessio Farcomeni ("Tor Vergata" University of Rome - Italy)

Luca Frigau (University of Cagliari - Italy)

Luis Ángel García Escudero (University of Valladolid - Spain)

Bettina Grün (Vienna University of Economics and Business - Austria)

Salvatore Ingrassia (University of Catania - Italy)

Volodymyr Melnykov (University of Alabama - USA)

Brendan Murphy (University College Dublin - Ireland)

Maria Lucia Parrella (University of Salerno - Italy)

Carla Rampichini (University of Florence - Italy)

Monia Ranalli (Sapienza University of Rome - Italy)

J. Sunil Rao (University of Miami - USA)

Marco Riani (University of di Parma - Italy)

Nicola Salvati (University of Pisa - Italy)

Laura Maria Sangalli (Polytechnic University of Milan - Italy)

Bruno Scarpa (University of Padua - Italy)

Mariangela Sciandra (University of Palermo - Italy)

Luca Scrucca (University of Perugia - Italy)

Domenico Vistocco (Federico II University of Naples - Italy)

Mariangela Zenga (University of Milan-Bicocca - Italy)

## LOCAL PROGRAM COMMITTEE

Carla Rampichini (chair) (University of Florence - Italy)

Chiara Bocci (University of Florence - Italy)

Anna Gottard (University of Florence - Italy)

Leonardo Grilli (University of Florence - Italy)

Monia Lupparelli (University of Florence - Italy)

Maria Francesca Marino (University of Florence - Italy)

Agnese Panzera (University of Florence - Italy)

Emilia Rocco (University of Florence - Italy)

Domenico Vistocco (Federico II University of Naples - Italy)

CLADAG 2021  
BOOK OF ABSTRACTS  
AND SHORT PAPERS

13th Scientific Meeting of the Classification  
and Data Analysis Group  
Firenze, September 9-11, 2021

edited by  
Giovanni C. Porzio  
Carla Rampichini  
Chiara Bocci

FIRENZE UNIVERSITY PRESS  
2021

CLADAG 2021 BOOK OF ABSTRACTS AND SHORT PAPERS : 13th Scientific Meeting of the Classification and Data Analysis Group Firenze, September 9-11, 2021/ edited by Giovanni C. Porzio, Carla Rampichini, Chiara Bocci. — Firenze : Firenze University Press, 2021.  
(Proceedings e report ; 128)

<https://www.fupress.com/isbn/9788855183406>

ISSN 2704-601X (print)

ISSN 2704-5846 (online)

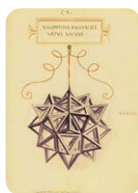
ISBN 978-88-5518-340-6 (PDF)

ISBN 978-88-5518-341-3 (XML)

DOI 10.36253/978-88-5518-340-6

Graphic design: Alberto Pizarro Fernández, Lettera Meccanica SRLs

Front cover: Illustration of the statue by Giambologna, *Appennino* (1579-1580) by Anna Gottard



Classification and Data  
Analysis Group (CLADAG)  
of the Italian Statistical  
Society (SIS)

*FUP Best Practice in Scholarly Publishing* (DOI [https://doi.org/10.36253/fup\\_best\\_practice](https://doi.org/10.36253/fup_best_practice))

All publications are submitted to an external refereeing process under the responsibility of the FUP Editorial Board and the Scientific Boards of the series. The works published are evaluated and approved by the Editorial Board of the publishing house, and must be compliant with the Peer review policy, the Open Access, Copyright and Licensing policy and the Publication Ethics and Complaint policy.

Firenze University Press Editorial Board

M. Garzaniti (Editor-in-Chief), M.E. Alberti, F. Vittorio Arrigoni, E. Castellani, F. Ciampi, D. D'Andrea, A. Dolfi, R. Ferrise, A. Lambertini, R. Lanfredini, D. Lippi, G. Mari, A. Mariani, P.M. Mariano, S. Marinai, R. Minuti, P. Nanni, A. Orlandi, I. Palchetti, A. Perulli, G. Pratesi, S. Scaramuzzi, I. Stolzi.

📄 The online digital edition is published in Open Access on [www.fupress.com](http://www.fupress.com).

Content license: except where otherwise noted, the present work is released under Creative Commons Attribution 4.0 International license (CC BY 4.0: <http://creativecommons.org/licenses/by/4.0/legalcode>). This license allows you to share any part of the work by any means and format, modify it for any purpose, including commercial, as long as appropriate credit is given to the author, any changes made to the work are indicated and a URL link is provided to the license.

Metadata license: all the metadata are released under the Public Domain Dedication license (CC0 1.0 Universal: <https://creativecommons.org/publicdomain/zero/1.0/legalcode>).

© 2021 Author(s)

Published by Firenze University Press  
Firenze University Press  
Università degli Studi di Firenze  
via Cittadella, 7, 50144 Firenze, Italy  
[www.fupress.com](http://www.fupress.com)

*This book is printed on acid-free paper  
Printed in Italy*

## INDEX

<b>Preface</b>	<b>1</b>
----------------	----------

### Keynote Speakers

<i>Jean-Michel Loubes</i> <b>Optimal transport methods for fairness in machine learning</b>	<b>5</b>
<i>Peter Rousseeuw, Jakob Raymaekers and Mia Hubert</i> <b>Class maps for visualizing classification results</b>	<b>6</b>
<i>Robert Tibshirani, Stephen Bates and Trevor Hastie</i> <b>Understanding cross-validation and prediction error</b>	<b>7</b>
<i>Cinzia Viroli</i> <b>Quantile-based classification</b>	<b>8</b>
<i>Bin Yu</i> <b>Veridical data science for responsible AI: characterizing V4 neurons through deepTune</b>	<b>9</b>

### Plenary Session

<i>Daniel Diaz</i> <b>A simple correction for COVID-19 sampling bias</b>	<b>14</b>
<i>Jeffrey S. Morris</i> <b>A seat at the table: the key role of biostatistics and data science in the COVID-19 pandemic</b>	<b>15</b>
<i>Bhramar Mukherjee</i> <b>Predictions, role of interventions and the crisis of virus in India: a data science call to arms</b>	<b>16</b>
<i>Danny Pfeffermann</i> <b>Contributions of Israel's CBS to rout COVID-19</b>	<b>17</b>

### Invited Papers

<i>Claudio Agostinelli, Giovanni Saraceno and Luca Greco</i> <b>Robust issues in estimating modes for multivariate torus data</b>	<b>21</b>
<i>Emanuele Aliverti</i> <b>Bayesian nonparametric dynamic modeling of psychological traits</b>	<b>25</b>

<i>Andres M. Alonso, Carolina Gamboa and Daniel Peña</i> <b>Clustering financial time series using generalized cross correlations</b>	27
<i>Raffaele Argiento, Edoardo Filippi-Mazzola and Lucia Paci</i> <b>Model-based clustering for categorical data via Hamming distance</b>	31
<i>Antonio Balzanella, Antonio Irpino and Francisco de A.T. De Carvalho</i> <b>Mining multiple time sequences through co-clustering algorithms for distributional data</b>	32
<i>Francesco Bartolucci, Fulvia Pennoni and Federico Cortese</i> <b>Hidden Markov and regime switching copula models for state allocation in multiple time-series</b>	36
<i>Michela Battauz and Paolo Vidoni</i> <b>Boosting multidimensional IRT models</b>	40
<i>Matteo Bottai</i> <b>Understanding and estimating conditional parametric quantile models</b>	44
<i>Niklas Bussmann, Roman Enzmann, Paolo Giudici and Emanuela Raffinetti</i> <b>Shapley Lorenz methods for eXplainable artificial intelligence</b>	45
<i>Andrea Cappelozzo, Ludovic Duponchel, Francesca Greselin and Brendan Murphy</i> <b>Robust classification of spectroscopic data in agri-food: first analysis on the stability of results</b>	49
<i>Andrea Cerasa, Enrico Checchi, Domenico Perrotta and Francesca Torti</i> <b>Issues in monitoring the EU trade of critical COVID-19 commodities</b>	53
<i>Marcello Chiodi</i> <b>Smoothed non linear PCA for multivariate data</b>	54
<i>Roberto Colombi, Sabrina Giordano and Maria Kateri</i> <b>Accounting for response behavior in longitudinal rating data</b>	58
<i>Claudio Conversano, Giulia Contu, Luca Frigau and Carmela Cappelli</i> <b>Network-based semi-supervised clustering of time series data</b>	62
<i>Federica Cugnata, Chiara Brombin, Pietro Cippà, Alessandro Ceschi, Paolo Ferrari and Clelia Di Serio</i> <b>Characterising longitudinal trajectories of COVID-19 biomarkers within a latent class framework</b>	64
<i>Silvia D'Angelo</i> <b>Sender and receiver effects in latent space models for multiplex data</b>	68
<i>Anna Denkowska and Stanisław Wanat</i> <b>DTW-based assessment of the predictive power of the copula-DCC-GARCH-MST model developed for European insurance institutions</b>	71
<i>Roberto Di Mari, Zsuzsa Bakk, Jennifer Oser and Jouni Kuha</i> <b>Two-step estimation of multilevel latent class models with covariates</b>	75
<i>Marie Du Roy de Chaumaray and Matthieu Marbac</i> <b>Clustering data with non-ignorable missingness using semi-parametric mixture models</b>	79

# HIDDEN MARKOV AND REGIME SWITCHING COPULA MODELS FOR STATE ALLOCATION IN MULTIPLE TIME-SERIES

Francesco Bartolucci<sup>1</sup>, Fulvia Pennoni<sup>2</sup>, and Federico P. Cortese<sup>3</sup>

<sup>1</sup> Department of Economics, University of Perugia  
(e-mail: francesco.bartolucci@unipg.it)

<sup>2</sup> Department of Statistics and Quantitative Methods, University of Milano-Bicocca  
(e-mail: fulvia.pennoni@unimib.it)

<sup>3</sup> Department of Economics, Management and Statistics, University of Milano-Bicocca  
(e-mail: f.cortese5@campus.unimib.it)

**ABSTRACT:** We consider hidden Markov and regime-switching copula models as approaches for state allocation in multiple time-series, where state allocation means prediction of the latent state characterizing each time occasion based on the observed data. This dynamic clustering, performed under the two model specifications, takes the correlation structure of the time-series into account. Maximum likelihood estimation of the model parameters is carried out by the expectation-maximization algorithm. For illustration we use data on the market of cryptocurrencies characterized by periods of high turbulence in which interdependence among assets is marked.

**KEYWORDS:** daily log-returns, expectation-maximization algorithm, forecast, latent variables, model-based clustering

## 1 Introduction

In the analysis of multiple time-series, state allocation, namely prediction of the state or regime underlying the observed data at a certain time occasion, is an important task, especially in finance and related fields. This type of clustering is dynamic because a different state may be predicted at every time occasion and may be based on models representing each time-specific state by a discrete latent variable assuming, typically, a few possible values. In this contribution, we compare two different model specifications of this type: multivariate hidden Markov (HM) models (Zucchini *et al.*, 2017) and regime-switching (RS) copulas (Rodriguez, 2007).

Among HM models we consider, in particular, those based on the assumption that the time-specific vector of observable variables follows a conditional Gaussian distribution with parameters depending on the latent state.



RS copulas are instead based on a copula function, which may be chosen among the Clayton, the Gumbel, the Gaussian, or the Student- $t$ , with parameters governed by a hidden Markov process of first-order so as to flexibly account for the correlation patterns between each pair of series.

The expectation-maximization (EM) algorithm (Dempster *et al.*, 1977) is used for maximum likelihood estimation of the parameters of both models. Model selection is performed to choose the most appropriate number of hidden states and evaluate the level of chain homogeneity over time (Bartolucci *et al.*, 2013). For the HM model, this selection is based on the Bayesian Information Criterion (BIC), and for RS copulas, it is also based on a goodness-of-fit procedure relying on the Cramér-von Mises statistic.

As an illustration we consider the problem of state allocation in analyzing time-series of the main cryptocurrencies daily log-returns over a three-year period.

## 2 Hidden Markov and Regime-Switching Copula Models

Let  $\mathbf{y}_t$ ,  $t = 1, 2, \dots$ , be the vector where each element  $y_{tj}$ ,  $j = 1, \dots, r$ , corresponds to the value of time-series  $j$  at time occasion  $t$ , with  $r$  denoting the number of time-series under consideration. The main assumption of the multivariate HM model is that the random vectors  $\mathbf{y}_1, \mathbf{y}_2, \dots$  are conditionally independent given a hidden process  $u_1, u_2, \dots$  that follows a first-order Markov chain with  $k$  states, labeled from 1 to  $k$ . This process is governed by the initial probabilities  $\pi_u = p(u_1 = u)$ ,  $u = 1, \dots, k$ , and the transition probabilities  $\pi_{u|\bar{u}} = p(u_t = u | u_{t-1} = \bar{u})$ ,  $t = 2, \dots$ ,  $\bar{u}, u = 1, \dots, k$ . We assume a Gaussian distribution for the observations at every time occasion, that is,  $\mathbf{y}_t | u_t = u \sim N_r(\boldsymbol{\mu}_u, \boldsymbol{\Sigma}_u)$ , where  $\boldsymbol{\mu}_u$  and  $\boldsymbol{\Sigma}_u$  are the mean vector and variance-covariance matrix for latent state  $u$ . The above assumptions imply that the conditional distribution of the time-series  $\mathbf{y}_1, \mathbf{y}_2, \dots$ , given the sequence of hidden states, may be expressed as  $f(\mathbf{y}_1, \mathbf{y}_2, \dots | u_1, u_2, \dots) = \prod_t \phi(\mathbf{y}_t; \boldsymbol{\mu}_{u_t}, \boldsymbol{\Sigma}_{u_t})$ , where  $\phi(\cdot; \cdot)$  denotes the density of the multivariate Gaussian distribution. The manifest distribution of the multiple time-series has the following density function:

$$f(\mathbf{y}_1, \mathbf{y}_2, \dots) = \sum_{u_1} \pi_{u_1} \phi(\mathbf{y}_1; \boldsymbol{\mu}_{u_1}, \boldsymbol{\Sigma}_{u_1}) \sum_{u_2} \pi_{u_2|u_1} \phi(\mathbf{y}_2; \boldsymbol{\mu}_{u_2}, \boldsymbol{\Sigma}_{u_2}) \cdots$$

Concerning the copula model, we first consider only the bivariate case, so we define  $\mathbf{y}_t = (y_{t1}, y_{t2})$  as a vector with elements  $y_{tj}$ ,  $j = 1, 2$ , corresponding to the observation for time-series  $j$  at time  $t = 1, 2, \dots$  and  $F_1$  and  $F_2$  as the

marginal cdfs of each time-series. Sklar’s theorem (Sklar, 1959) allows us to separate the fitting of the marginal cdfs from the fitting of the joint distribution, represented by a copula function. This approach consists in estimating the two marginal distributions, obtaining  $\hat{F}_1$  and  $\hat{F}_2$ , and then computing the normalized ranks of the pseudo-observations  $\tilde{\mathbf{e}}_t = (\tilde{e}_{t1}, \tilde{e}_{t2})$  as  $\tilde{e}_{tj} = \text{rank}(\hat{z}_{tj}) / (T + 1)$ , with  $\hat{z}_{tj} = \hat{F}_j(y_{tj})$ , and  $T$  being the number of observed time occasions. Finally, for the pseudo-observations  $\tilde{\mathbf{e}}_t$ , an RS copula model is assumed based on a hidden homogeneous Markov process denoted as  $v_1, v_2, \dots$ , with  $k$  states. The copula density indicated with  $c(\cdot; \cdot)$  may be chosen among the Clayton, the Gumbel, the Gaussian, or the Student- $t$  copulas, with state-specific parameter  $\beta_v$ . The density of the pseudo-observations is given by

$$f(\tilde{\mathbf{e}}_1, \tilde{\mathbf{e}}_2, \dots) = \sum_{v_1} \pi_{v_1} c(\tilde{\mathbf{e}}_1; \beta_{v_1}) \sum_{v_2} \pi_{v_2|v_1} c(\tilde{\mathbf{e}}_2; \beta_{v_2}) \cdots,$$

and it is based on the initial and transition probabilities defined as above.

Given that the state sequence is not observable, a full maximum likelihood approach for estimating the parameters of both models is carried out through the EM algorithm. Following the current literature, model selection for the HM model is based on the BIC, and for the RS copula it is also performed through a goodness-of-fit procedure consisting in calculating a  $p$ -value referred to the Cramér-von Mises statistic for the hypothesis of correct model specification.

We compare the performance of HM models and RS copulas focusing on the crucial aspect of state allocation. The optimal state allocation is performed by finding the optimal joint sequence  $\tilde{u}_1, \tilde{u}_2, \dots$  (or  $\tilde{v}_1, \tilde{v}_2, \dots$ ) of unknown states given the corresponding observations. This clustering procedure, also known as global decoding, is achieved through the Viterbi algorithm (Viterbi, 1967), which is a dynamic programming algorithm.

We also aim at extending the RS copula approach to an arbitrary number of time-series  $r$  rather than to only 2. In this regard, we propose the composite likelihood approach (Varin *et al.*, 2011) for estimation, which is based on considering all possible ordered pairs of time-series among the available ones.

### 3 Application

As an illustration, for the HM model we consider the joint daily log-returns\* of the five cryptocurrencies Bitcoin, Ethereum, Ripple, Litecoin, and Bitcoin

\*provided by the Crypto Asset Lab: <https://cryptoassetlab.diseade.unimib.it/>.

Cash, for the period 2017-2020. For the RS copulas, allowing only for bivariate associations, we define four copulas where the bivariate vector of observations consists of the Bitcoin and each of the other four cryptocurrencies. Results for the HM model show that the minimum value of the BIC is reached considering a five-state heteroschedastic structure. According to these estimates, there are three negative regimes (in terms of estimated expected log-returns), with relatively high and positive correlations of Bitcoin with all the other cryptocurrencies, and two states with positive returns and lower correlations. Regarding the global decoding, these two states are the most likely in the first year of observation, and the other three states characterize the last two years.

Concerning the RS copulas, and considering as an example the couple of cryptocurrencies Bitcoin-Ethereum, we observe that a three-regime Clayton copula provides the best fit. Given that the Clayton copula allows for explicit computation of the lower tail correlation index, we estimate that two regimes provide zero or low values for the lower tail index, and the third regime provides high values for it. According to the optimal state sequence, we estimate that there is substantial interchangeability between the first two states in the whole period, whereas the third state is the most likely for the last year of observation.

## References

- BARTOLUCCI, F., FARCOMENI, A., & PENNONI, F. 2013. *Latent Markov Models for Longitudinal Data*. Boca Raton, FL: Chapman & Hall/CRC.
- DEMPSTER, A. P., LAIRD, N. M., & RUBIN, D. B. 1977. Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society, Series B*, **39**, 1–38.
- RODRIGUEZ, J. C. 2007. Measuring financial contagion: A copula approach. *Journal of Empirical Finance*, **14**, 401–423.
- SKLAR, M. 1959. Fonctions de repartition à n dimensions et leurs marges. *Publications de l'Institut Statistique de l'Université de Paris*, **8**, 229–231.
- VARIN, C., REID, N., & FIRTH, D. 2011. An overview of composite likelihood methods. *Statistica Sinica*, **21**, 5–42.
- VITERBI, A. 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory*, **13**, 260–269.
- ZUCCHINI, W., MACDONALD, I. L., & LANGROCK, R. 2017. *Hidden Markov Models for Time Series: An Introduction Using R*. Boca Raton, FL: CRC.