

PROGRAMME AND ABSTRACTS

CFE-CMStatistics 2025

19th International Conference on
Computational and Financial Econometrics (CFE 2025)

and

Computational and Methodological Statistics (CMStatistics 2025)

<https://www.cmstatistics.org/CFECMStatistics2025>

Birkbeck, University of London, UK

13–15 December 2025



ISBN 978-9925-7812-9-4

©2025 - ECOSTA ECONOMETRICS AND STATISTICS

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted, in any other form or by any means without the prior permission from the publisher.

Co-chairs CFE:

Esther Ruiz, Michael Smith and Weining Wang.

Co-chairs CMStatistics:

Joshua Cape, Sonja Greven, Matias Quiroz and Mark Steel.

International Organizing Committee:

Ana Colubi, George Kapetanios, Erricos Kontoghiorghes and David Weston.

Scientific Programme Committee CFE:

Alessandra Amendola, Jonas Andersson, Christian Brownlees, Andrew Butters, Massimiliano Caporin, Angeles Carnero, Gianluca Cubadda, Christina Erlwein-Sayer, Domenico Giannone, Luigi Grossi, Roxana Halbleib, Chen Huang, Joann Jasiak, Lingwei Kong, Robinson Kruse-Becher, Carlos Lamarche, Nathan Lassurance, Michele Lenza, Johan Lyhagen, Malvina Marchese, Kamiar Mohaddes, Jose Olmo, Michael Owyang, Dario Palumbo, Peter Pedroni, Markus Pelger, Juan Manuel Rodriguez Poo, Yongcheol Shin, Genaro Sucarrat, Hideatsu Tsukahara, Ralf Wilke and Yohei Yamamoto.

Scientific Programme Committee CMStatistics:

Eleonora Arnone, Boris Beranger, Alexandros Beskos, Matteo Borrotti, Fan Bu, Nilanjan Chakraborty, Hao Chen, Shanshan Ding, Matteo Fontana, Debarghya Ghoshdastidar, Steven Gilmour, Bettina Gruen, Arthur Guillaumin, Mayetri Gupta, Vanda Inacio, Galin Jones, Cai Li, Tianxi Li, Po-Ling Loh, Vince Lyzinski, Etienne Marceau, Wendy Meiring, Ramses Mena, Silvia Montagna, Anirbit Mukherjee, Alejandro Murua, Anahita Nodehi, Lucia Paci, Simone Padoan, Alessia Pini, Yixuan Qiu, Marialuisa Restaino, Nicholas Rios, Sanjoy Sinha, Tiejun Tong, Minh-Ngoc Tran, Masayuki Uchida, Chenguang Wang, Wei-Ying Wu and Ruqing Zhu.

Local Organizers:

Birkbeck, University of London, CFEnetwork and CMStatistics.

Dear Friends and Colleagues,

We warmly welcome you to London for the 19th International Conference on Computational and Financial Econometrics (CFE 2025) and Computational and Methodological Statistics (CMStatistics 2025).

The conference aims to bring together researchers and practitioners to discuss recent methodology and computational approaches for economics, finance, and statistics. The CFE-CMStatistics 2025 programme consists of about 315 sessions, four plenary talks, and nearly 1300 presentations. With about 1400 participants, this conference stands out as one of the most important international scientific events in the field.

CFE 2025 hosts an expanded section honoring Prof. Hashem Pesaran. This includes his keynote talk and over forty contributions from several of his notable collaborators in a series of special commemorative invited sessions.

The co-chairs have endeavoured to provide a balanced and stimulating programme that will appeal to the diverse interests of the participants. The international organizing committee hopes that the hybrid conference will provide an ideal environment to communicate effectively with colleagues. The conference is the collective effort of many individuals and organizations. The Scientific Programme Committee, the Session Organizers, the supporting universities, and various agents have contributed substantially to the organization of the conference. We acknowledge their work and the support of our networks.

Birkbeck, University of London offers excellent facilities and a fantastic environment in central London. Through their efforts, the local hosts and sponsoring organizations have substantially contributed to the successful organization of the conference. We thank them all for their support. In particular, we express our sincere appreciation to the host, the School of Computing and Mathematical Sciences, Birkbeck, University of London.

We are pleased to announce that the 2024 impact factor of the official journal of CFEnetwork and CMStatistics, *Econometrics and Statistics (EcoSta)*, stands at 2.5. *EcoSta* is ranked Q1 in both Economics and in Statistics and Probability. CMStatistics also publishes *The Annals of Statistical Data Science (SDS)* as a supplement to the Elsevier journal *Computational Statistics & Data Analysis (CSDA)*. The CSDA is the official journal of CMStatistics as well. CSDA continues to uphold its commendable and consistent performance, with an impact factor of 1.6. You are encouraged to submit your papers to *EcoSta*, the *Annals of SDS* or regular peer-reviewed issues of CSDA.

Looking ahead, the CFE-CMStatistics 2026 conference will be hosted at HTW Berlin, University of Applied Sciences, from Saturday, December 12th, to Monday, December 14th, 2026, with tutorials scheduled prior to the conference. We extend a heartfelt invitation and enthusiastic encouragement for your active participation in these forthcoming events.

We wish you a productive and stimulating conference.

Kind regards,

Ana Colubi, Erricos J. Kontoghiorghes, and David Weston
International Organizing Committee

**CMStatistics: ERCIM Working Group on
COMPUTATIONAL AND METHODOLOGICAL STATISTICS**

<http://www.cmstatistics.org>

The working group (WG) CMStatistics comprises a number of specialized teams in various research areas of computational and methodological statistics. The teams act autonomously within the framework of the WG in order to promote their own research agenda. Their activities are endorsed by the WG. They submit research proposals, organize sessions, tracks and tutorials during the annual WG meetings and edit journal special issues. The Econometrics and Statistics (EcoSta) and Computational Statistics & Data Analysis (CSDA) are the official journals of the CMStatistics.

Specialized teams

Currently, the ERCIM WG has over 1950 members and the following specialized teams

| | |
|--|--|
| BIO: Biostatistics | NPS: Non-Parametric Statistics |
| BS: Bayesian Statistics | RS: Robust Statistics |
| DMC: Dependence Models and Copulas | SA: Survival Analysis |
| DOE: Design Of Experiments | SAE: Small Area Estimation |
| FDA: Functional Data Analysis | SDS: Statistical Data Science: Methods and Computations |
| HDS: High-Dimensional Statistics | SEA: Statistics of Extremes and Applications |
| IS: Imprecision in Statistics | SL: Statistical Learning |
| LVSEM: Latent Variable and Structural Equation Models | TSMC: Times Series |
| MM: Mixture Models | |

You are encouraged to become a member of the WG. For further information, please contact the Chairs of the specialized groups (see the WG's website) or email at info@cmstatistics.org.

**CFEnetwork
COMPUTATIONAL AND FINANCIAL ECONOMETRICS**

<http://www.CFEnetwork.org>

The Computational and Financial Econometrics (CFEnetwork) comprises a number of specialized teams in various research areas of theoretical and applied econometrics, financial econometrics and computation, and empirical finance. The teams contribute to the network's activities by organizing sessions, tracks and tutorials during the annual CFEnetwork meetings, and by submitting research proposals. Furthermore, the teams edit special issues currently published under the Annals of CFE. The Econometrics and Statistics (EcoSta) is the official journal of the CFEnetwork. Currently, the CFEnetwork has over 1100 members.

You are encouraged to become a member of the CFEnetwork. For further information, please see the website or contact by email at info@cfenetwork.org.

SCHEDULE (UTC+0)

| 2025-12-13 | 2025-12-14 | 2025-12-15 |
|--|---|--|
| Opening , 08:30 - 08:45 | | |
| A - Keynote CFE-CMStatistics 2025 08:45 - 09:35 | F CFE-CMStatistics 2025 08:45 - 10:25 | K - Keynote CFE-CMStatistics 2025 09:10 - 10:00 |
| Coffee break 09:35 - 10:05 | Coffee break 10:25 - 10:55 | Coffee break 10:00 - 10:30 |
| B CFE-CMStatistics 2025 10:05 - 12:10 | G CFE-CMStatistics 2025 10:55 - 12:10 | L CFE-CMStatistics 2025 10:30 - 12:10 |
| Lunch break 12:10 - 13:40 | Lunch break 12:10 - 13:40 | Lunch break 12:10 - 13:40 |
| C CFE-CMStatistics 2025 13:40 - 15:20 | H CFE-CMStatistics 2025 13:40 - 15:20 | M - Keynote CFE-CMStatistics 2025 13:40 - 14:30 |
| D - Keynote CFE-CMStatistics 2025 15:30 - 16:20 | Coffee break 15:20 - 15:50 | N CFE-CMStatistics 2025 14:40 - 16:20 |
| Coffee break 16:20 - 16:50 | I CFE-CMStatistics 2025 15:50 - 17:30 | Coffee break 16:20 - 16:50 |
| E CFE-CMStatistics 2025 16:50 - 18:55 | J CFE-CMStatistics 2025 17:40 - 18:55 | O CFE-CMStatistics 2025 16:50 - 18:30 |
| Welcome reception 19:00 - 20:30 | | Closing networking drink 18:45 - 19:45 |
| | Conference buffet dinner and DJ party 20:00 - 23:00 | |

TUTORIALS, MEETINGS AND CONFERENCE DETAILS (see maps)

TUTORIALS

Four independent tutorials will take place from the 10th to the 12th of December 2025, organized within the framework of the COST Action HiTEc. The first tutorial, “Probabilistic programming for statistical analysis in Julia”, will be coordinated by Prof. Mattias Villani, Stockholm University, Sweden. The second tutorial, “Bayesian variable selection”, will be coordinated by Prof. Jim Griffin, University College London, UK. The third tutorial, “Small ball probabilities for functional data analysis”, will be coordinated by Prof. Enea Bongiorno, Università del Piemonte Orientale, Italy. The fourth tutorial, “Measure transportation, statistical inference, and time series”, will be coordinated by Prof. Marc Hallin, Université libre de Bruxelles, Belgium. Further details are available on the website. Only participants who have subscribed to the tutorials can attend, either in person or virtually through the conference website.

ECOSTA and CSDA EDITORIAL BOARD MEETINGS

The *Econometrics and Statistics (EcoSta) Editorial Board* and the *CSDA and Annals of Statistical Data Science Editorial Board* meetings will be held from 18:45 to 20:00 on Monday, 15 December, at the first floor of the Fitzroy Tavern, 16 Charlotte St., London W1T 2LY.

CONFERENCE DETAILS

Access to the conference and social events

- Participants can attend virtually or in person according to what they selected when registering.
- The in-person venue is Birkbeck, University of London, comprising Birkbeck Central Building (BCB) and the Birkbeck Malet Street (main building). Please note that Birkbeck Central opens at 8:00, so you will not be able to enter before that time.
- All the parallel sessions will take place at Birkbeck Central Building (BCB) and the keynote talks will take place at Birkbeck Malet Street (main building). Participants can use BCB 305 Upper Hall as quiet room to attend sessions virtually with their laptops and headphones.
- The **registration** desk will be located at the entrance hall of Birkbeck Central. Registration will be open on Friday afternoon, from 16:00 to 18:00, during the weekend from 8:00 to 18:00 and on Monday from 8:30 to 16:30.
- The **coffee breaks** will take place in the Conservatory located on the first floor of Birkbeck Central.
- The **welcome reception** will take place in the Conservatory located on the first floor of Birkbeck Central.
- The **conference dinner** will take place at The Memoir Club (formerly Ambassadors Bloomsbury), 12 Upper Woburn Pl, London WC1H 0HX.
- The **closing networking drink** will take place at the first floor of Fitzroy Tavern, 16 Charlotte St., London W1T 2LY.
- For environmental sustainability reasons, the conference aims to minimize paper usage and overall consumption. While there will be a limited number of printed Books of Abstracts, bags, pens, and pads available for those who requested them during their registration, we strongly encourage all participants to opt for digital materials by downloading them onto their personal devices. QR codes will be displayed in the registration area. These codes will enable participants to quickly access essential information, further reducing the need for printed materials and promoting a paperless conference experience.
- The conference is live-streamed, and it will not be recorded. The presentations will take place through Zoom. The conference programme time is set at UTC+0.
- In order to access the virtual conference, you must first log in to the registration tool, get the daily password there, and leave the session open. Then you should open another tab and go to the interactive programme (schedule). Click on the slot you wish to attend and then on the session. You will be redirected to Zoom, where you will need to use the daily password.
- Please note that for security reasons, the Zoom links will not be sent to the speakers, and they can only be found on the online interactive programme (schedule).
- Detailed instructions for virtual and in-person attendance, hybrid sessions, speakers, chairs, posters, networking, test sessions, as well as FAQs, can be found on the webpage.

Presentation instructions

The paper presentations must be shared through Zoom. The in-person rooms will be visible in Zoom as the corresponding hybrid session. Virtual speakers should have a stable internet connection, and make sure their video and audio are working. In-person participants can bring presentations on a USB in PDF format or download them from personal clouds/emails on the room's PC. Alternatively, speakers can use their own laptop to connect to the corresponding Zoom session. They will share their slides when the chair requires it, present their talk, and be ready to answer questions after the presentation. Detailed instructions for speakers can be found on the website. Standard talks have 20 minutes for the talk and 3-4 minutes for discussion as a general rule. Strict timing must be observed.

Posters

Posters will be displayed in a common Zoom session and do not need to be uploaded or sent in advance. In-person participants can join the poster session from the room designated in the programme with their devices. The Chair will give time to the presenters to show their posters, which will be rotated sequentially to be seen by the session audience. Presenters may use their time to briefly talk about their posters or simply show them and wait for questions or comments from the audience. Detailed instructions for the poster presentations can be found on the website.

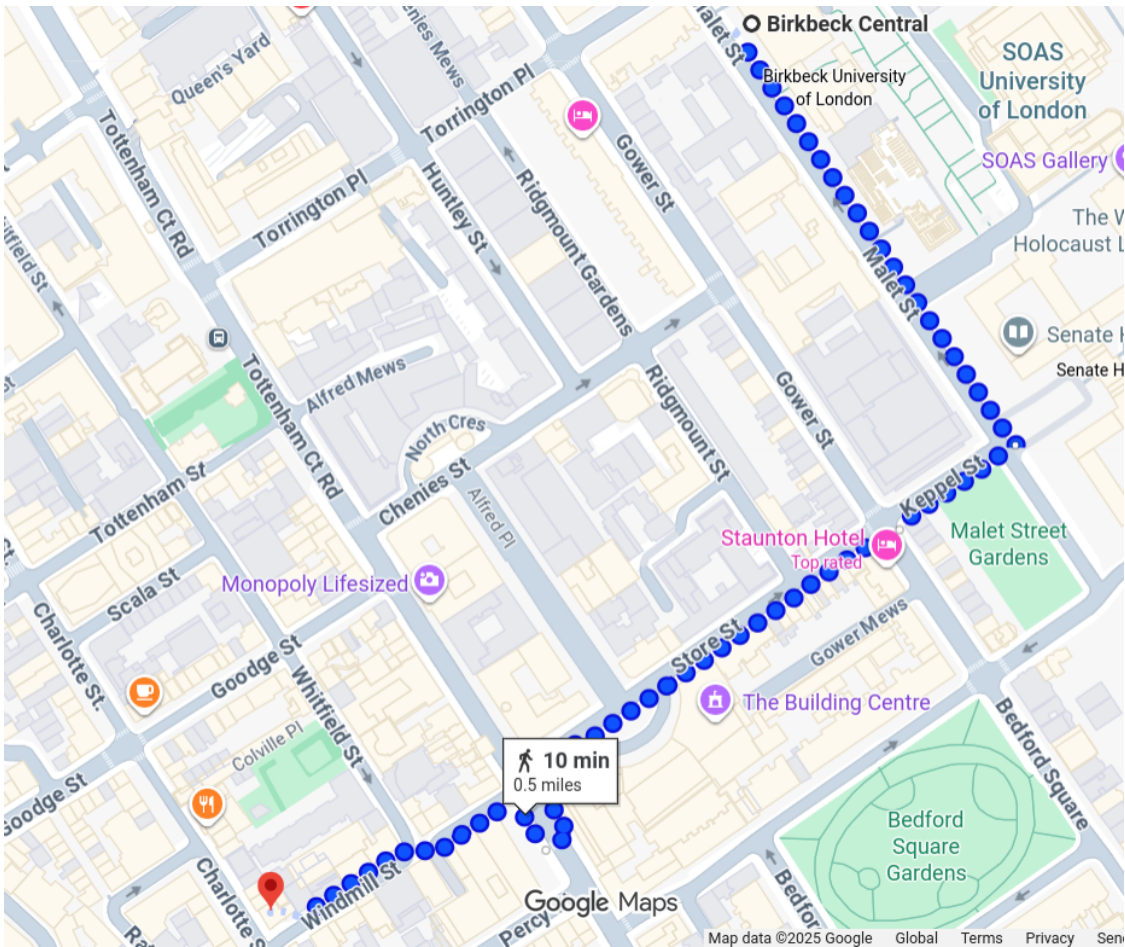
Session chairs

Session chairs will be responsible for introducing the session and speakers, and coordinating discussion time. A conference staff member, identified on Zoom as “Angel”, will assist online. In case of a missing speaker, or technical problem with a speaker, the Chair can move to the next speaker and return later if possible. Detailed instructions for session chairs in both virtual and hybrid sessions can be found on the website.

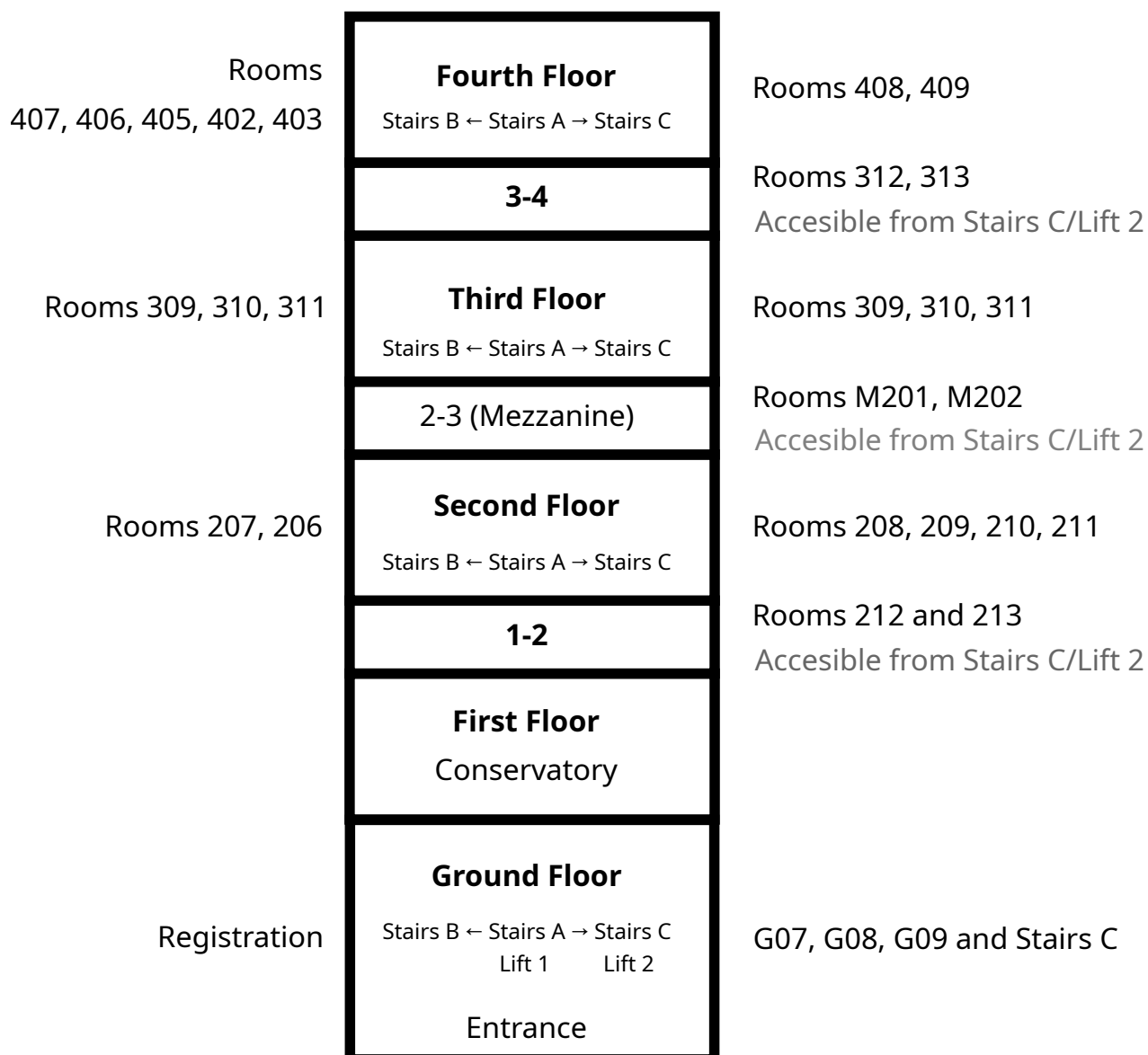
Conference locations map



Map for the closing drink



Vertical map of Birkbeck Central



PUBLICATION OUTLETS

Econometrics and Statistics (EcoSta)

<http://www.elsevier.com/locate/ecosta>

Econometrics and Statistics (EcoSta), published by Elsevier, is the official journal of the networks Computational and Financial Econometrics and Computational and Methodological Statistics. It publishes research papers on all aspects of econometrics and statistics and comprises two sections: **Part A: Econometrics.** Emphasis is given to methodological and theoretical papers containing substantial econometrics derivations or showing potential for a significant impact in the broad area of econometrics. Topics of interest include the estimation of econometric models and associated inference, model selection, panel data, measurement error, Bayesian methods, and time series analyses. Simulations are considered when they involve an original methodology. Innovative papers in financial econometrics and its applications are considered. The covered topics include portfolio allocation, option pricing, quantitative risk management, systemic risk and market microstructure. Interest is focused as well on well-founded applied econometric studies that demonstrate the practicality of new procedures and models. Such studies should involve the rigorous application of statistical techniques, including estimation, inference and forecasting. Topics include volatility and risk, credit risk, pricing models, portfolio management, and emerging markets. Innovative contributions in empirical finance and financial data analysis that use advanced statistical methods are encouraged. The results of the submissions should be replicable. Applications consisting only of routine calculations are not of interest to the journal.

Part B: Statistics. Papers providing important original contributions to methodological statistics inspired in applications are considered for this section. Papers dealing, directly or indirectly, with computational and technical elements are particularly encouraged. These cover developments concerning issues of high-dimensionality, re-sampling, dependence, robustness, filtering, and, in general, the interaction of mathematical methods, numerical implementations and the extra burden of analysing large and/or complex datasets with such methods in different areas such as medicine, epidemiology, biology, psychology, climatology and communication. Innovative algorithmic developments are also of interest, as are the computer programs and the computational environments that implement them as a complement.

The journal consists, preponderantly, of original research. Occasionally, reviews and short papers from experts are published, which may be accompanied by discussions. Special issues and sections within important areas of research are occasionally published. The journal publishes as a supplement the Annals of Computational and Financial Econometrics.

Call For Papers Econometrics and Statistics (EcoSta)

<http://www.elsevier.com/locate/ecosta>

Papers presented at the conference and containing novel components in econometrics or statistics are encouraged to be submitted for publication in special peer-reviewed or regular issues of the Elsevier journal Econometrics and Statistics (EcoSta) and its supplement Annals of Computational and Financial Econometrics. Papers should be submitted using the EM Submission tool. In the EM please select as type of article the CFE conference, CMStatistics Conference or Annals of Computational and Financial Econometrics. Any questions may be directed via email to editor@econometricsandstatistics.org

Call For Papers CSDA Annals of Statistical Data Science (SDS)

<http://www.elsevier.com/locate/csda>

We are inviting submissions for the 1st issue of the CSDA Annals of Statistical Data Science. The Annals of Statistical Data Science is published as a supplement to the journal of Computational Statistics & Data Analysis. It will serve as an outlet for distinguished research papers using advanced computational and/or statistical methods for tackling challenging data analytic problems. The Annals will become a valuable resource for well-founded theoretical and applied data-driven research. Authors submitting a paper to CSDA may request that it be considered for inclusion in the Annals. Each issue will be assigned to several Guest Associate Editors who will be responsible, together with the CSDA Co-Editors, for the selection of papers.

Submissions for the Annals should contain a significant computational or statistical methodological component for data analytics. In particular, the Annals welcomes contributions at the interface of computing and statistics, addressing problems involving large and/or complex data. Emphasis will be given to comprehensive and reproducible research, including data-driven methodology, algorithms and software. There is no deadline for submissions. Papers can be submitted at any time. When they have been received, they will enter the editorial system immediately. All submissions must contain original unpublished work not being considered for publication elsewhere. Please submit your paper electronically using the Elsevier Editorial System: <http://ees.elsevier.com/csda> (Choose Article Type: Research paper, and then Select "Section IV. Annals of Statistical Data Science").

Contents

| | |
|--|--------------------------------------|
| General Information | I |
| Committees | III |
| Welcome | IV |
| CMStatistics: ERCIM Working Group on Computational and Methodological Statistics | V |
| CFEnetwork: Computational and Financial Econometrics | V |
| Scientific programme | V |
| Tutorials, Meetings and Conference details | VII |
| Conference locations and closing drink maps | VIII |
| Vertical map of Birkbeck Central | IX |
| Publications outlets of the journals EcoSta and CSDA and Call for papers | X |
| Keynote Talks | 1 |
| Keynote talk I (Marc Hallin, Universite Libre de Bruxelles, Belgium) | Saturday 13.12.2025 at 08:45 - 09:35 |
| Multiple-attribute Lorenz functions and Gini indices: A measure-transportation approach | 1 |
| Keynote talk II (Esther Ruiz, Universidad Carlos III de Madrid, Spain) | Saturday 13.12.2025 at 15:30 - 16:20 |
| Scenarios for macroeconomic risk | 1 |
| Keynote talk IV (Hashem Pesaran, Trinity College, Cambridge, United Kingdom) | Monday 15.12.2025 at 09:10 - 10:00 |
| The output convergence debate revisited: Lessons from recent developments in the analysis of panel data models | 1 |
| Keynote talk III (Jim Griffin, University College London, United Kingdom) | Monday 15.12.2025 at 13:40 - 14:30 |
| Normalized latent measure models | 1 |
| Parallel Sessions | 2 |
| Parallel Session B – CFE-CMStatistics 2025 (Saturday 13.12.2025 at 10:05 - 12:10) | 2 |
| CI015: DENSITY DATA ANALYSIS (Room: BCB 307) | 2 |
| CO154: HiTEC: MACHINE LEARNING IN ECONOMETRICS (Room: BCB G07) | 2 |
| CO137: MACROECONOMIC RISK (Room: BCB G08) | 3 |
| CO094: ADVANCES IN MACRO- AND FINANCIAL ECONOMETRICS (Room: BCB G09) | 3 |
| CO296: ADVANCES IN PANEL DATA AND TIME SERIES MODELS FOR FINANCIAL ECONOMETRICS (Room: Virtual R01) | 4 |
| CO036: CFE SESSION: A TRIBUTE TO H. PESARAN I (Room: BCB 206) | 5 |
| CO219: STATISTICAL METHODS FOR INSURANCE (Room: BCB 207) | 5 |
| CO266: MODERN STATISTICAL METHODS FOR THE ANALYSIS OF COMPLEX REAL DATA (Room: BCB 208) | 6 |
| CO208: INFERENCE AND CHANGE DETECTION FOR DEPENDENT DATA (Room: BCB 209) | 7 |
| CO133: STATISTICS FOR STOCHASTIC PROCESSES AND THEIR APPLICATIONS (Room: BCB 210) | 7 |
| CO121: BAYESIAN ECONOMETRICS (Room: BCB 211) | 8 |
| CO180: COMPUTATIONAL METHODS IN STATISTICS (Room: BCB 212) | 9 |
| CO044: ADVANCES IN QUANTITATIVE RISK MANAGEMENT AND COPULAS (Room: BCB 213) | 10 |
| CO293: ASYMPTOTICS, STATISTICS OF EXTREMES AND QUANTILE INFERENCE (Room: BCB M201) | 10 |
| CO415: CLIMATE FINANCE (Room: BCB M202) | 11 |
| CO153: STATISTICAL MODELING OF TEXT DATA (Room: BCB 308) | 12 |
| CO185: HIGH-DIMENSIONAL DATA ANALYSIS (Room: BCB 309) | 13 |
| CO025: STATISTICAL METHODS AND THEIR APPLICATIONS (Room: BCB 310) | 13 |
| CO141: COMPLEX DATA UNCOVERED: METHODOLOGY AND PRACTICAL INSIGHTS (Room: BCB 311) | 14 |
| CO053: LATENT STRUCTURE IN COMPLEX DATA: BAYESIAN AND FREQUENTIST VIEWS (Room: BCB 403) | 15 |
| CO239: MODEL-BASED CLUSTERING: THEORY AND APPLICATIONS (Room: BCB 405) | 16 |
| CO332: ADVANCES IN MODERN CAUSAL INFERENCE (Room: BCB 407) | 17 |
| CO342: SPATIAL STATISTICS, IMAGE ANALYSIS AND DIRECTIONAL DATA (Room: BCB 409) | 18 |
| CC350: TIME SERIES INFERENCE AND ESTIMATION (Room: BCB 312) | 18 |
| CC357: SURVIVAL ANALYSIS (Room: BCB 313) | 19 |
| CC367: BIostatISTICS (Room: BCB 402) | 20 |
| CC363: NONPARAMETRIC METHODS (Room: BCB 406) | 21 |
| CC358: RISK ANALYSIS (Room: BCB 408) | 21 |
| Parallel Session C – CFE-CMStatistics 2025 (Saturday 13.12.2025 at 13:40 - 15:20) | 23 |
| CI017: ECOSta PART B: STATISTICS INVITED SESSION (VIRTUAL) (Room: BCB 307) | 23 |
| CO186: HiTEC: SOME NEW CHALLENGES IN FUNCTIONAL STATISTICS (Room: BCB G07) | 23 |
| CO236: RECENT CONTRIBUTIONS TO RISK ANALYSIS (Room: BCB G08) | 24 |
| CO032: ML AND HIGH-DIMENSIONAL MACROECONOMIC FORECASTING MODELS (Room: BCB G09) | 24 |
| CO218: ADVANCED METHODS IN SINGLE-CELL GENOMICS (Room: Virtual R01) | 25 |
| CO054: RECENT ADVANCES IN ANALYZING NETWORK DATA (VIRTUAL) (Room: Virtual R02) | 25 |
| CO019: CFE SESSION: A TRIBUTE TO H. PESARAN (Room: BCB 206) | 26 |
| CO177: TRACKING BUSINESS CYCLES IN HIGH DIMENSIONAL AND/OR VOLATILE SETTINGS (Room: BCB 207) | 26 |
| CO416: NON-LINEAR IMPULSES AND PROPAGATION MECHANISMS (Room: BCB 208) | 27 |
| CO155: ADVANCES IN MACHINE LEARNING FOR COMMODITY FORECASTING (Room: BCB 209) | 27 |
| CO035: RECENT DEVELOPMENTS IN LONG-RUN MODELLING (Room: BCB 210) | 28 |
| CO331: MODELLING TIME-VARYING RELATIONSHIPS IN ECONOMETRICS (Room: BCB 211) | 29 |

| | |
|---|-----------|
| CO024: NON-GAUSSIAN TIME SERIES (Room: BCB 212) | 29 |
| CO101: RECENT ADVANCES IN QUANTILE REGRESSION (Room: BCB 213) | 30 |
| CO038: STATISTICAL INFERENCE AND PREDICTION OF EXTREME VALUES AND BEYOND (Room: BCB M201) | 30 |
| CO192: THE FINANCIAL RISKS OF CLIMATE CHANGE AND BIODIVERSITY LOSS (Room: BCB M202) | 31 |
| CO090: STATISTICAL LEARNING WITH MIXTURES (Room: BCB 308) | 31 |
| CO245: STARSTRUCK STATISTICS (Room: BCB 309) | 32 |
| CO083: ADVANCES IN LONGITUDINAL AND REPEATED MEASURES DATA (Room: BCB 310) | 33 |
| CO163: RECENT ADVANCES IN THE DESIGN OF EXPERIMENTS (Room: BCB 312) | 33 |
| CO136: MODERN DIMENSION REDUCTION TECHNIQUES AND THEORIES (Room: BCB 313) | 34 |
| CO291: STATISTICAL METHODS FOR LARGE-SCALE BIOMEDICAL DATA ANALYSIS (Room: BCB 402) | 34 |
| CO105: COMPLEX PROBLEMS IN CAUSAL INFERENCE (Room: BCB 403) | 35 |
| CO248: RECENT DEVELOPMENTS IN BAYESIAN CLUSTERING (Room: BCB 405) | 36 |
| CO104: RECENT ADVANCEMENTS IN STATISTICAL NETWORK ANALYSIS (Room: BCB 406) | 37 |
| CO157: CAUSAL INFERENCE AND PERSONALIZED MEDICINE (Room: BCB 407) | 37 |
| CO303: ADVANCES IN BAYESIAN METHODS AND COMPUTATIONS (Room: BCB 408) | 38 |
| CO202: DEVELOPMENTS IN SPATIO-TEMPORAL DISEASE MAPPING AND SURVEILLANCE (Room: BCB 409) | 38 |
| CC444: SHORT TALKS: CMSTATISTICS I (Room: BCB 311) | 39 |
| Parallel Session E – CFE-CMStatistics 2025 (Saturday 13.12.2025 at 16:50 - 18:55) | 41 |
| CI246: CSDA STATISTICAL DATA SCIENCE (Room: BCB 307) | 41 |
| CO168: ADVANCES IN STATISTICAL MACHINE LEARNING FOR MODERN DATA CHALLENGES (Room: BCB G07) | 41 |
| CO193: BAYESIAN TIME-SERIES MODELLING (Room: BCB G08) | 42 |
| CO075: TOPICS IN FINANCIAL ECONOMETRICS (Room: BCB G09) | 43 |
| CO139: STATISTICAL ADVANCES IN MACHINE LEARNING FOR COMPLEX BIOMEDICAL DATA (Room: Virtual R01) | 44 |
| CO278: CFE SESSION: A TRIBUTE TO H. PESARAN (Room: BCB 206) | 44 |
| CO232: MACHINE LEARNING METHODS FOR NOWCASTING (Room: BCB 207) | 45 |
| CO241: NONCAUSAL TIME SERIES (Room: BCB 209) | 45 |
| CO171: ADVANCES IN TIME SERIES FOR ECONOMICS AND FINANCE (Room: BCB 210) | 46 |
| CO018: ECOSta JOURNAL SESSION (Room: BCB 211) | 47 |
| CO050: ADVANCES IN RISK MEASUREMENT AND IN FORECASTING (Room: BCB 212) | 48 |
| CO255: HIGH-DIMENSIONAL AND STRUCTURAL APPROACHES TO SYSTEMIC RISK ASSESSMENT (Room: BCB 213) | 48 |
| CO194: NEW TOOLS FOR CAUSAL INFERENCE AND FORECASTING (Room: BCB M202) | 49 |
| CO100: RESEARCHER DEGREES OF FREEDOM, FLEXIBLE MODELING AND INTERPRETABILITY (Room: BCB 308) | 50 |
| CO095: TOPICS IN STATISTICAL GENETICS AND BIOINFORMATICS (Room: BCB 309) | 51 |
| CO065: NEW DEVELOPMENTS IN FUNCTIONAL DATA ANALYSIS (Room: BCB 310) | 52 |
| CO264: STATISTICAL METHODS FOR ECOLOGICAL AND ENVIRONMENTAL COMPLEXITY (Room: BCB 311) | 52 |
| CO327: STATISTICAL AI IDEAS FOR MODERN BIOSTATISTICS (Room: BCB 312) | 53 |
| CO112: DEVELOPMENTS IN RECEIVER OPERATING CHARACTERISTIC CURVE ANALYSIS (Room: BCB 313) | 54 |
| CO110: SELECTED WORK FROM CAUCUS FOR WOMEN IN STATISTICS (CWS) MEMBERS (Room: BCB 402) | 55 |
| CO112: CAUSAL INFERENCE FOR POLICY EVALUATION AND TREATMENT DECISION MAKING (Room: BCB 403) | 56 |
| CO056: BAYESIAN NONPARAMETRIC MIXTURE MODELING AND CLUSTERING (VIRTUAL) (Room: BCB 405) | 56 |
| CO020: ADVANCES IN MODELLING AND INFERENCE ON NETWORKS AND HYPERGRAPHS (Room: BCB 406) | 57 |
| CO067: ADVANCES IN CAUSAL INFERENCE AND MEDIATION WITH COMPLEX DATA (Room: BCB 407) | 58 |
| CO148: STATISTICAL METHODS FOR BIOMEDICAL APPLICATIONS (Room: BCB 409) | 59 |
| CC425: STATISTICAL COMPUTING AND SIMULATION (Room: BCB 208) | 59 |
| CC439: ASSET PRICING AND TESTING (Room: BCB M201) | 60 |
| CC404: APPLIED STATISTICS (Room: BCB 408) | 61 |
| Parallel Session F – CFE-CMStatistics 2025 (Sunday 14.12.2025 at 08:45 - 10:25) | 62 |
| CI011: ADVANCES IN BAYESIAN COMPUTATION FOR COMPLEX MODELS (Room: BCB 307) | 62 |
| CO052: HiTEC: ADVANCES IN ECONOMETRIC METHODS FOR ECONOMICS AND FINANCE (Room: BCB G07) | 62 |
| CO107: APPLICATIONS IN HIGH-DIMENSIONAL DATASETS (Room: BCB G08) | 63 |
| CO073: MACROECONOMIC POLICIES AND MACROECONOMETRICS (Room: BCB G09) | 63 |
| CO029: BAYESIAN NONPARAMETRICS FOR TIME SERIES ANALYSIS (Room: Virtual R01) | 64 |
| CO077: ARDL MODELS AND COINTEGRATION: A TRIBUTE TO H. PESARAN (Room: BCB 206) | 64 |
| CO172: STATISTICAL MODELS FOR SHADOW ECONOMY AND FRAUD DETECTION (Room: BCB 209) | 65 |
| CO049: CFE SESSION: A TRIBUTE TO H. PESARAN III (Room: BCB 211) | 66 |
| CO134: ADVANCES IN TIME SERIES MODELING (Room: BCB 212) | 66 |
| CO084: THEORY AND COMPUTATION FOR STOCHASTIC PROCESS MODELS (Room: BCB 213) | 67 |
| CO238: ROBUST METHODS IN ENERGY, ENVIRONMENT, AND FINANCE (Room: BCB M202) | 68 |
| CO027: STATISTICAL LEARNING FOR COMPLEX AND DYNAMIC SYSTEMS (Room: BCB 308) | 68 |
| CO138: MODERN STATISTICAL LEARNING AND INFERENCE FOR STRUCTURED DATA (Room: BCB 309) | 69 |
| CO037: ADVANCES IN FUNCTIONAL AND COMPLEX DATA ANALYSIS (Room: BCB 310) | 69 |
| CO267: CLUSTERING OF HETEROGENEOUS DATA (Room: BCB 311) | 70 |
| CO126: DESIGN OF EXPERIMENTS METHODOLOGY AND APPLICATIONS (Room: BCB 312) | 70 |
| CO174: SURVIVAL ANALYSIS METHODS FOR PARTLY INTERVAL-CENSORED DATA (Room: BCB 313) | 71 |
| CO275: SPATIO-TEMPORAL MODELLING FOR ENVIRONMENTAL AND CLIMATE APPLICATIONS (Room: BCB 402) | 72 |

| | |
|---|-----------|
| CO085: ADVANCES IN STATISTICS AND ML FOR BREEDING STUDIES (Room: BCB 403) | 72 |
| CO320: TOPICS IN BAYESIAN MODELLING AND COMPUTATION (Room: BCB 405) | 73 |
| CO310: SCALABLE INFERENCE IN LARGE NETWORKS (Room: BCB 406) | 74 |
| CO210: SPORTS ANALYTICS (Room: BCB 408) | 74 |
| CC426: MACROECONOMIC FORECASTING AND UNCERTAINTY (Room: BCB 207) | 75 |
| CC437: FINANCIAL ECONOMETRICS I (Room: BCB 210) | 76 |
| CC360: ROBUST METHODS (Room: BCB M201) | 76 |
| CC430: STATISTICAL ALGORITHMS AND MACHINE LEARNING TECHNIQUES (Room: BCB 407) | 77 |
| CC419: SPATIAL STATISTICS (Room: BCB 409) | 77 |
| CP001: POSTER SESSION (Room: BCB 208) | 78 |
| Parallel Session G – CFE-CMStatistics 2025 (Sunday 14.12.2025 at 10:55 - 12:10) | 81 |
| CO028: HiTEC: STATISTICAL INFERENCE ON INSTABILITIES AND INDEPENDENCE (Room: BCB G07) | 81 |
| CO188: PHYSICS-GUIDED SPATIOTEMPORAL DATA: DESCRIPTIVE VS DYNAMICAL SYSTEMS (Room: BCB G08) | 81 |
| CO069: FINANCIAL AND MACRO ECONOMETRIC PREDICTABILITY (Room: BCB G09) | 81 |
| CO234: ADVANCEMENT ON STATISTICAL NETWORK MODELING AND ANALYSIS (Room: Virtual R01) | 82 |
| CO114: CFE SESSION: A TRIBUTE TO H. PESARAN IV (Room: BCB 206) | 82 |
| CO196: ECONOMETRIC FORECASTING (Room: BCB 207) | 83 |
| CO191: MODELLING RISK AND UNCERTAINTY (Room: BCB 209) | 83 |
| CO340: ADVANCES IN NON-LINEAR TIME SERIES (Room: BCB 210) | 83 |
| CO295: STATISTICAL MODELS FOR BUSINESS AND ECONOMICS (Room: BCB 211) | 84 |
| CO127: ECONOMIC UNCERTAINTY AND ITS EFFECTS (Room: BCB 212) | 85 |
| CO338: THE ECONOMETRICS OF EMPIRICAL ASSET PRICING (Room: BCB 213) | 85 |
| CO106: EXTREMES AND DEPENDENCE MODELLING (Room: BCB M201) | 85 |
| CO184: BAYESIAN METHODS IN FINANCE (Room: BCB M202) | 86 |
| CO041: OVER-PARAMETRIZATION AND OVERFITTING IN MACHINE LEARNING (Room: BCB 308) | 86 |
| CO206: SAFE AI FOR DECISION-MAKING IN ECONOMICS AND FINANCE (VIRTUAL) (Room: BCB 309) | 87 |
| CO315: INTERACTING URNS AND INNOVATION PROCESSES (Room: BCB 310) | 87 |
| CO276: UNIFORM INFERENCE ON HIGH-DIMENSIONAL MODELS (Room: BCB 311) | 88 |
| CO322: STATISTICAL METHODS FOR ENVIRONMENTAL SUSTAINABILITY (Room: BCB 312) | 88 |
| CO198: AI FOR COMPLEX NETWORKS: APPLICATIONS IN HEALTH AND FINANCIAL SYSTEMS (Room: BCB 403) | 89 |
| CO328: RECENT ADVANCES OF REINFORCEMENT LEARNING AND DYNAMIC DECISION-MAKING (Room: BCB 407) | 89 |
| CO181: J-ISBA SESSION: NEW ADVANCES IN BAYESIAN STATISTICS (Room: BCB 408) | 90 |
| CC420: ESG SCORES AND SUSTAINABLE FINANCE (Room: BCB 208) | 90 |
| CC433: HIGH-DIMENSIONAL AND SPARSE MODELING (Room: BCB 307) | 91 |
| CC407: EMPIRICAL FINANCE (Room: BCB 313) | 91 |
| CC431: APPLIED FINANCIAL ECONOMETRICS (Room: BCB 402) | 92 |
| CC436: APPLIED STATISTICAL TOOLS AND EVALUATION (Room: BCB 405) | 92 |
| CC408: METHODOLOGICAL STATISTICS (Room: BCB 406) | 93 |
| CC429: HIGH-FREQUENCY AND STOCHASTIC ECONOMETRICS (Room: BCB 409) | 93 |
| Parallel Session H – CFE-CMStatistics 2025 (Sunday 14.12.2025 at 13:40 - 15:20) | 95 |
| CI010: DYNAMIC FACTOR MODELS AND BEYOND (Room: BCB 211) | 95 |
| CO068: HiTEC: ADVANCES IN FINANCIAL ECONOMETRICS (Room: BCB G07) | 95 |
| CO059: APPLIED MACRO (Room: BCB G08) | 96 |
| CO058: STATISTICAL AND ECONOMETRIC MODELS FOR CENSORED DATA (Room: BCB G09) | 96 |
| CO304: RECENT ADVANCES IN CAUSAL INFERENCE (Room: Virtual R01) | 97 |
| CO286: STATISTICAL MODELING AND INFERENCE UNDER COMPLEX STRUCTURE (VIRTUAL) (Room: Virtual R02) | 97 |
| CO080: NONLINEARITIES AND APPLICATIONS: A TRIBUTE TO H. PESARAN (VIRTUAL) (Room: BCB 206) | 98 |
| CO021: INNOVATIONS IN FINANCE AND INSURANCE (Room: BCB 207) | 99 |
| CO250: PRICES AND WAGES INFLATION (Room: BCB 208) | 99 |
| CO169: DEPENDENCE MODELING IN ACTUARIAL SCIENCE (Room: BCB 210) | 100 |
| CO306: ADVANCES IN ASSET PRICING (Room: BCB 212) | 101 |
| CO339: ASYMPTOTICS IN INFERENCE AND COMPUTATION FOR STOCHASTIC PROCESSES (Room: BCB 213) | 101 |
| CO277: MULTIVARIATE EXTREMES (Room: BCB M201) | 102 |
| CO298: COMPUTATIONAL METHODS IN FINANCIAL ECONOMICS (Room: BCB M202) | 102 |
| CO197: DATA DEPTH IN STATISTICS AND MACHINE LEARNING (Room: BCB 307) | 103 |
| CO314: STATISTICAL CHANGE-POINT DETECTION (Room: BCB 310) | 103 |
| CO156: ADVANCES IN OPTIMAL DESIGN AND SUBSAMPLING (VIRTUAL) (Room: BCB 312) | 104 |
| CO299: MODERN APPROACHES TO DISTRIBUTIONAL AND HIGH-DIMENSIONAL DATA ANALYSIS (Room: BCB 403) | 104 |
| CO323: ADVANCES IN BAYESIAN GRAPHICAL MODEL INFERENCE (Room: BCB 405) | 105 |
| CO116: RECENT ADVANCES ON NETWORK AND MATRIX DATA ANALYSIS (Room: BCB 406) | 106 |
| CO076: BAYESIAN SEMI- AND NON-PARAMETRIC METHODS II (Room: BCB 408) | 106 |
| CO108: STATISTICS AND SPORT (Room: BCB 409) | 107 |
| CC391: TEXT DATA IN ECONOMICS AND FINANCE (Room: BCB 209) | 108 |
| CC379: MACHINE LEARNING (Room: BCB 308) | 108 |
| CC443: SHORT TALKS: CFE (Room: BCB 309) | 109 |

| | |
|--|------------|
| CC418: CLUSTERING (Room: BCB 311) | 110 |
| CC428: COMPLEX AND STRUCTURED TIME SERIES (Room: BCB 313) | 111 |
| CC384: ALGORITHMS AND SOFTWARE (Room: BCB 402) | 111 |
| CC423: NONPARAMETRIC HYPOTHESIS TESTS (Room: BCB 407) | 112 |
| Parallel Session I – CFE-CMStatistics 2025 (Sunday 14.12.2025 at 15:50 - 17:30) | 113 |
| CO103: EMBEDDINGS IN STATISTICS AND AI (Room: BCB 307) | 113 |
| CO151: HiTEC: STATISTICAL LEARNING AND REGIME-SWITCHING IN FINANCE (Room: BCB G07) | 113 |
| CO317: ADVANCES IN STRUCTURAL ECONOMETRICS (Room: BCB G08) | 113 |
| CO099: TOPICS IN MACROFINANCE AND ECONOMETRICS (Room: BCB G09) | 114 |
| CO269: EMPIRICAL MACRO (Room: Virtual R01) | 115 |
| CO042: CFE SESSION: A TRIBUTE TO H. PESARAN II (Room: BCB 206) | 115 |
| CO132: SL FOR FORECASTING, STRUCTURAL REDUCTION, AND REGIME CHANGE (Room: BCB 207) | 116 |
| CO122: ADVANCES IN COMPOSITE LIKELIHOOD INFERENCE FOR HIGH-DIMENSIONAL DATA (Room: BCB 208) | 117 |
| CO066: DIRECTIONAL STATISTICS (Room: BCB 209) | 117 |
| CO093: TIME SERIES AND STRUCTURAL BREAKS (Room: BCB 210) | 118 |
| CO022: LEARNING AND DISCOVERY IN HIGH-DIMENSIONAL AND STRUCTURED DATA (Room: BCB 211) | 118 |
| CO289: TREATMENT EFFECTS AND POLICY EVALUATION (Room: BCB 212) | 119 |
| CO294: ADVANCES IN EXTREME VALUE ANALYSIS (Room: BCB M201) | 119 |
| CO082: SUSTAINABLE AND CLIMATE FINANCE (Room: BCB M202) | 120 |
| CO119: THE PRICE OF INSIGHT: COST-AWARE AND HUMAN-GUIDED LEARNING (Room: BCB 308) | 121 |
| CO251: MARKOV-SWITCHING MODELS (Room: BCB 309) | 121 |
| CO124: RECENT ADVANCES IN LEARNING AND COMPLEX BIOMEDICAL DATA (Room: BCB 310) | 122 |
| CO087: EMERGING TOPICS IN STATISTICAL META-ANALYSIS AND APPLICATIONS (Room: BCB 311) | 122 |
| CO051: DESIGN AND ANALYSIS OF EXPERIMENTS (Room: BCB 312) | 123 |
| CO209: BAYESIAN METHODS IN ENVIRONMENTAL, EPIDEMIOLOGIC, GENOMIC PROBLEMS (Room: BCB 402) | 124 |
| CO225: ADVANCES IN STATISTICAL ANALYSES OF NETWORK DATA (Room: BCB 403) | 124 |
| CO070: MONTE CARLO METHODS AND THEIR APPLICATION (Room: BCB 405) | 125 |
| CO103: NETWORK ANALYSIS METHODS AND THEIR APPLICATIONS (Room: BCB 406) | 126 |
| CO285: INNOVATIONS IN CAUSAL, PREDICTIVE, AND NONPARAMETRIC INFERENCE (Room: BCB 407) | 126 |
| CO118: BAYESIAN SEMI- AND NON-PARAMETRIC METHODS I (Room: BCB 408) | 127 |
| CO091: STATISTICAL METHODS AND APPLICATION FOR HIGH DIMENSIONAL BIOMARKERS (Room: BCB 409) | 128 |
| CC442: SHORT TALKS: CMSTATISTICS II (Room: BCB 213) | 128 |
| CC400: STATISTICAL MODELLING (Room: BCB 313) | 129 |
| CV394: FORECASTING (VIRTUAL) (Room: Virtual R02) | 130 |
| Parallel Session J – CFE-CMStatistics 2025 (Sunday 14.12.2025 at 17:40 - 18:55) | 131 |
| CO064: BAYESIAN MIXTURE AND LATENT VARIABLE MODELS FOR LARGE-SCALE PROBLEMS (Room: BCB G07) | 131 |
| CO089: RECENT ADVANCES IN FUNCTIONAL DATA ANALYSIS (Room: Virtual R01) | 131 |
| CO023: NEW MACHINE LEARNING AND BAYESIAN TECHNIQUES (VIRTUAL) (Room: Virtual R02) | 132 |
| CO222: DYNAMIC PANEL DATA MODELS: A TRIBUTE TO H. PESARAN (Room: BCB 206) | 132 |
| CO147: NONPARAMETRIC METHODS FOR DEPENDENCE AND STRUCTURAL CHANGE (Room: BCB 207) | 132 |
| CO123: QUANTIFYING MODEL SELECTION UNCERTAINTY (Room: BCB 208) | 133 |
| CO131: IMPUTATION TECHNIQUES AND SPATIO-TEMPORAL MODELING (Room: BCB 209) | 133 |
| CO047: SPATIO-TEMPORAL DEPENDENCE AND COPULA MODELS (Room: BCB 210) | 134 |
| CO062: ADVANCES IN MULTIVARIATE TIME SERIES (Room: BCB 211) | 134 |
| CO111: ADVANCES IN TIME SERIES ANALYSIS AND FORECASTING (Room: BCB 213) | 135 |
| CO233: TOPICS IN FINANCE AND ECONOMETRICS (Room: BCB M202) | 135 |
| CO162: RECENT ADVANCES IN HIDDEN MARKOV MODELS (Room: BCB 307) | 136 |
| CO173: STATISTICAL INNOVATIONS FOR DEPENDENT DATA (Room: BCB 308) | 136 |
| CO316: RECENT ADVANCES IN MATRIX AND TENSOR TIME SERIES MODELS (Room: BCB 309) | 136 |
| CO125: RECENT APPROACHES TO ENVIRONMENTAL AND SPATIO-TEMPORAL STATISTICS (Room: BCB 310) | 137 |
| CO201: BRIDGING CAUSALITY, CONNECTIVITY, AND ROBUST INFERENCE (Room: BCB 311) | 138 |
| CO200: INFERENCE AND TESTING IN COMPLEX BIOMEDICAL STUDIES (Room: BCB 312) | 138 |
| CO257: RECENT ADVANCES IN STATISTICAL LEARNING ON COMPLEX DATA SETS (Room: BCB 313) | 139 |
| CO135: MODERN STATISTICAL FRONTIERS FOR BIOMEDICAL AND HIGH-DIMENSIONAL DATA (Room: BCB 402) | 139 |
| CO150: STATISTICAL MODELING FOR BIOMEDICAL AND PUBLIC HEALTH APPLICATIONS (Room: BCB 403) | 140 |
| CO179: HIGH-DIMENSIONAL AND DYNAMIC BAYES: MODELS, COMPUTATION, APPLICATIONS (Room: BCB 405) | 140 |
| CO337: NOVEL METHODS FOR INFERRING NETWORK STRUCTURE (Room: BCB 406) | 141 |
| CO166: BAYESIAN NONPARAMETRIC METHODS FOR LARGE-SCALE INFERENCE PROBLEMS (Room: BCB 407) | 141 |
| CO143: STATISTICS IN NEUROSCIENCE II (Room: BCB 409) | 142 |
| CC435: MACROECONOMICS AND SOCIAL POLICY (Room: BCB G09) | 142 |
| CC424: CAUSAL INFERENCE AND POLICY EVALUATION (Room: BCB 212) | 143 |
| CC396: STOCHASTIC PROCESSES (Room: BCB M201) | 143 |
| CC440: APPLIED STATISTICAL ANALYSIS IN ECONOMICS AND SOCIETY (Room: BCB 408) | 144 |
| CV366: BIOSTATISTICS (VIRTUAL) (Room: BCB G08) | 144 |

| | |
|---|------------|
| Parallel Session L – CFE-CMStatistics 2025 (Monday 15.12.2025 at 10:30 - 12:10) | 145 |
| CI007: CFE SPECIAL INVITED SESSION: A TRIBUTE TO H. PESARAN II (Room: BCB 206) | 145 |
| CO140: HiTEC: RECENT ADVANCES IN MODEL SPECIFICATION TESTING (Room: BCB G07) | 145 |
| CO165: ADVANCED STATISTICAL METHODS FOR ENERGY ECONOMICS (Room: BCB G08) | 146 |
| CO220: CLEVER MODELS FOR COMPLICATED DATA (Room: BCB G09) | 146 |
| CO333: CONTEMPORARY ISSUES IN ECONOMICS AND FINANCE (Room: BCB 207) | 147 |
| CO324: INNOVATIONS IN BAYESIAN NETWORK PSYCHOMETRICS (Room: BCB 208) | 148 |
| CO319: CROSS-SECTIONAL ASSET PRICING (Room: BCB 210) | 148 |
| CO072: TIME SERIES ECONOMETRICS (Room: BCB 211) | 149 |
| CO061: DYNAMIC MODELS FOR FINANCIAL DATA ANALYSES (Room: BCB 212) | 149 |
| CO043: RECENT ADVANCES IN ASYMPTOTIC STATISTICS FOR STOCHASTIC PROCESSES (Room: BCB 213) | 150 |
| CO235: NEW CONTRIBUTIONS TO EXTREME VALUE THEORY (Room: BCB M201) | 150 |
| CO086: MACHINE LEARNING AND HIGH-DIMENSIONAL DATA (Room: BCB 308) | 151 |
| CO231: TOPOLOGICAL AND GEOMETRIC STATISTICAL ANALYSIS (Room: BCB 309) | 152 |
| CO211: RECENT ADVANCES IN SURVIVAL DATA ANALYSIS (Room: BCB 313) | 152 |
| CO274: VARIABLE IMPORTANCE AND STATISTICAL LEARNING IN ENVIRONMENT AND ENERGY (Room: BCB 402) | 153 |
| CO161: BAYESIAN STRUCTURE LEARNING IN GRAPHICAL MODELS (Room: BCB 405) | 153 |
| CO273: BAYESIAN ANALYSIS OF NETWORK DATA (Room: BCB 406) | 154 |
| CO115: BAYESIAN SEMI- AND NON-PARAMETRIC METHODS III (Room: BCB 408) | 155 |
| CO074: SPATIO-TEMPORAL MODELING AND INFERENCE (Room: BCB 409) | 155 |
| CC371: PORTFOLIO OPTIMIZATION (Room: BCB 209) | 156 |
| CC421: TAIL-RISK MODELING (Room: BCB M202) | 156 |
| CC427: MACRO-FINANCIAL MODELING (Room: BCB 307) | 157 |
| CC375: FUNCTIONAL DATA ANALYSIS (Room: BCB 310) | 158 |
| CC432: FISCAL AND FINANCIAL ECONOMETRICS (Room: BCB 311) | 158 |
| CC422: ELECTRICITY MARKETS (Room: BCB 312) | 159 |
| CC388: CATEGORICAL DATA (Room: BCB 403) | 159 |
| CC387: CAUSAL INFERENCE AND NETWORK DATA (Room: BCB 407) | 160 |
| CV434: STATISTICAL MODELS AND INFERENCE (Room: Virtual R01) | 160 |
| Parallel Session N – CFE-CMStatistics 2025 (Monday 15.12.2025 at 14:40 - 16:20) | 162 |
| CI006: CFE SPECIAL INVITED SESSION: A TRIBUTE TO H. PESARAN I (Room: BCB 206) | 162 |
| CI016: CAUSAL INFERENCE UNDER UNOBSERVED CONFOUNDING (Room: BCB 307) | 162 |
| CO081: HiTEC: CLUSTERING OF COMPLEX DATA STRUCTURES (Room: BCB G07) | 162 |
| CO182: EMPIRICAL MACRO (Room: BCB G08) | 163 |
| CO120: RISK ANALYSIS AND MODEL SPECIFICATION (Room: BCB G09) | 164 |
| CO271: STATISTICAL MODELING AND MACHINE LEARNING FOR COMPLEX DATA (Room: Virtual R01) | 164 |
| CO252: ADVANCES IN FISCAL POLICY (Room: BCB 207) | 165 |
| CO292: ECONOMETRIC METHODS IN ENERGY, CLIMATE, AND RESOURCE RESEARCH (Room: BCB 208) | 165 |
| CO096: RECENT ADVANCES IN STATISTICAL FEW-SHOT LEARNING (Room: BCB 209) | 166 |
| CO301: NONCAUSAL ECONOMETRICS (Room: BCB 211) | 167 |
| CO321: HIGH-DIMENSIONAL AND NONPARAMETRIC METHODS FOR PANEL DATA MODELS (Room: BCB 212) | 167 |
| CO144: ADVANCES IN PROBABILISTIC FORECASTING (Room: BCB 213) | 168 |
| CO031: RECENT ADVANCES IN STATISTICAL FINANCE (Room: BCB M202) | 168 |
| CO290: ADVANCES IN STATISTICAL GENOMICS (Room: BCB 308) | 169 |
| CO199: RECENT ADVANCES IN MODEL SELECTION (Room: BCB 310) | 170 |
| CO224: INNOVATIONS IN STATISTICAL METHODS FOR COMPLEX DEPENDENCE (Room: BCB 311) | 170 |
| CO262: STATISTICAL RESEARCH IN DESIGN OF EXPERIMENTS (Room: BCB 312) | 171 |
| CO078: EMERGING TOPICS IN STATISTICS AND DATA SCIENCE (Room: BCB 313) | 172 |
| CO243: RECENT METHODOLOGICAL ADVANCES IN BIostatISTICS (Room: BCB 402) | 172 |
| CO092: STATISTICAL ANALYSES OF COMPLEX AND MULTIPLEX/MULTILAYER NETWORKS (Room: BCB 403) | 173 |
| CO308: STATISTICAL METHODS IN NETWORKS AND HIGH-DIMENSIONAL DATA (Room: BCB 406) | 173 |
| CO312: STATISTICAL INFERENCE AND LEARNING IN COMPLEX SYSTEMS (Room: BCB 407) | 174 |
| CO045: ADVANCES IN MONTE CARLO METHODS FOR BAYESIAN STATISTICS (Room: BCB 408) | 174 |
| CO187: RECENT ADVANCES IN CAUSAL ANALYSIS AND SPATIAL STATISTICS (Room: BCB 409) | 175 |
| CC399: ECONOMETRIC AND FINANCIAL MODELLING (Room: BCB 210) | 176 |
| CC438: FINANCIAL ECONOMETRICS II (Room: BCB M201) | 176 |
| CC376: APPLIED MACHINE LEARNING (Room: BCB 309) | 177 |
| CC347: BAYESIAN METHODS (Room: BCB 405) | 178 |

| | |
|--|------------|
| Parallel Session O – CFE-CMStatistics 2025 (Monday 15.12.2025 at 16:50 - 18:30) | 179 |
| CI009: CFE SPECIAL INVITED SESSION: A TRIBUTE TO H. PESARAN III (Room: BCB 206) | 179 |
| CO071: HiTEC: ADVANCES IN MATRIX TIME SERIES IN ECONOMETRICS (Room: BCB G07) | 179 |
| CO183: EMPIRICAL MACRO-FINANCE (Room: BCB G08) | 180 |
| CO057: TOPICS IN ECONOMETRICS (Room: BCB G09) | 180 |
| CO229: BAYESIAN METHODS FOR COMPOSITIONAL AND PROPORTIONAL DATA (Room: Virtual R01) | 181 |
| CO240: SYNTHETIC DATA: GENERATION AND VALIDATION METHODS (Room: BCB 208) | 181 |
| CO330: LARGE DATA ANALYSIS IN CANCER GENOMICS AND STATISTICS FOR CANCER (Room: BCB 209) | 182 |
| CO334: NETWORK ECONOMETRICS (Room: BCB 212) | 183 |
| CO178: RECENT ADVANCES IN PANEL TIME SERIES METHODS AND APPLICATIONS (Room: BCB 213) | 183 |
| CO414: INNOVATIVE STATISTICAL APPLICATIONS (Room: BCB M201) | 184 |
| CO128: MODELING, FORECASTING, AND POLICY ASSESSMENT: MACRO, FINANCE (Room: BCB M202) | 185 |
| CO302: ADVANCES IN STATISTICAL MACHINE LEARNING (Room: BCB 307) | 185 |
| CO272: APPLICATION OF MACHINE LEARNING IN SAMPLE SURVEYS AND SMALL AREA ESTIMATION (Room: BCB 308) | 186 |
| CO030: ADVANCES IN HIGH-DIMENSIONAL TIME SERIES AND BAYESIAN MODELING (Room: BCB 309) | 186 |
| CO326: STATISTICAL METHODOLOGY FOR DATA WITH SPATIAL AND TEMPORAL DEPENDENCIES (Room: BCB 310) | 187 |
| CO113: ROBUSTNESS AND REGULARIZATION IN MIXTURE MODELING AND NETWORK ANALYSIS (Room: BCB 311) | 188 |
| CO336: STATISTICAL GENETICS AND MULTIOMIC DATA ANALYSIS (Room: BCB 312) | 188 |
| CO325: ADVANCES IN SEQUENTIAL DECISION MAKING (Room: BCB 313) | 189 |
| CO158: NOVEL STATISTICAL AND COMPUTATIONAL METHODS FOR BIOMEDICAL SCIENCES (Room: BCB 402) | 190 |
| CO281: ADVANCED STATISTICAL METHODS FOR SPATIAL TRANSCRIPTOMICS DATA ANALYSIS (Room: BCB 403) | 190 |
| CO205: BAYESIAN STATISTICS IN BIOLOGICAL AND ENVIRONMENTAL SCIENCES (Room: BCB 405) | 191 |
| CO026: ADVANCES IN NETWORK ANALYSIS AND NONPARAMETRIC STATISTICS (Room: BCB 406) | 192 |
| CO318: THEORY AND APPLICATION OF SAMPLING ALGORITHMS (Room: BCB 408) | 192 |
| CO142: STATISTICS IN NEUROSCIENCE I (Room: BCB 409) | 193 |
| CC395: FORECASTING (Room: BCB 210) | 193 |
| CC403: APPLIED ECONOMETRICS (Room: BCB 211) | 194 |
| CC441: STATISTICAL INFERENCE AND COMPUTATION (Room: BCB 407) | 195 |
| CV413: FINANCIAL ECONOMETRICS (VIRTUAL) (Room: BCB 207) | 195 |

Saturday 13.12.2025 08:45 - 09:35

Room: MAL B34 Chair: Manfred Deistler

Keynote talk I

Multiple-attribute Lorenz functions and Gini indices: A measure-transportation approachSpeaker: **Marc Hallin, Université Libre de Bruxelles, Belgium**

Gilles Mordant

Based on the measure-transportation ideas and the related concepts of quantile functions and regions developed by a prior study, multiple-output generalizations of the traditional univariate concepts of Lorenz and concentration functions and related Gini and Kakwani coefficients are proposed. These new concepts have a natural interpretation, either in terms of contributions of quantile regions to the expectation of some variable of interest, or in terms of the physical notions of work and energy, which sheds new light on the nature of economic and social inequalities. When based on center-outward quantile regions, the proposed concepts pave the way to a statistically sound definition, based on multiple variables, of the notion, so far limited to bivariate characterizations, of middle-class, a notion of practical importance in various socio-economic and political contexts.

Saturday 13.12.2025 15:30 - 16:20

Room: MAL B34 Chair: Tommaso Proietti

Keynote talk II

Scenarios for macroeconomic riskSpeaker: **Esther Ruiz, Universidad Carlos III de Madrid, Spain**

The interest in developing econometric tools to estimate densities of key macroeconomic variables is clear. Point forecasts are insufficient for informed decision-making. Densities measure macroeconomic vulnerability, crucial for resilience policies. Recently, econometricians and policy-makers have been interested in constructing realistic scenarios that can help understand economic resilience by providing early warnings of adverse, low-probability but potentially catastrophic events. The aim is to propose estimating conditional densities under stressed-factor scenarios using factor-augmented quantile regressions (FA-QR) with factors extracted from dynamic factor models (DFMs) via principal components (PCs). Severe yet plausible stress scenarios are based on the joint distribution of these factors. The results are illustrated by calculating growth-in-stress (GiS) for US growth and the 5% quantile of stressed growth densities, and show GiS as a useful complement to growth-at-risk (GaR) for scenario analysis. The COVID-19 shock provides a natural environment to assess US growth vulnerability. Scenarios for US inflation are also obtained when domestic, international, and temperature factors are stressed, showing that, depending on the estimation method of factors' uncertainty, worst inflation-at-risk (IaR) scenarios can differ, with important implications for monetary policy design.

Monday 15.12.2025 09:10 - 10:00

Room: MAL B34 Chair: Ron Smith

Keynote talk IV

The output convergence debate revisited: Lessons from recent developments in the analysis of panel data modelsSpeaker: **Hashem Pesaran, Trinity College, Cambridge, United Kingdom**

Ron Smith

There has been a resurgence of interest in the debate on whether output per capita has been converging across countries, and the possible determinants of the cross-country differences in output per capita. The analysis of dynamic heterogeneous panels with interactive effects provides a fruitful framework within which the assumptions and econometric techniques employed in this literature can be examined. Much of the convergence debate assumes that country-specific outputs follow a partial adjustment process towards an unobserved steady state, partly driven by a latent factor taken as a proxy for global technology. Empirical studies within this and similar literatures tend to use panel or cross-section estimators that impose strong homogeneity restrictions on adjustment speeds and time effects, the so-called parallel trends assumption. Heterogeneity across units that is correlated with the regressors causes the estimates to be inconsistent, and such correlated heterogeneity is inherent in dynamic models. A systematic, step-by-step, theoretical analysis is conducted of the implications of parameter heterogeneity in intercepts, slopes, and factor loadings, showing how neglected heterogeneity can result in inconsistent estimates and distorted inference. It also examines the treatment of both time-varying and time-invariant covariates that could explain the differences in cross-country outputs. The theoretical results are illustrated using Penn World Table data.

Monday 15.12.2025 13:40 - 14:30

Room: MAL B34 Chair: Mark Steel

Keynote talk III

Normalized latent measure modelsSpeaker: **Jim Griffin, University College London, United Kingdom**

Normalized latent measure models (NLMMs) are a framework for modeling and comparing probability distributions using mixtures of nonparametric distributions. Although the methods are generic, the focus is on modeling similar distributions. For example, to model a large collection of probability distributions (such as areal income distributions) or the effects of covariates on a probability distribution (the effect of wind direction and speed on energy generation). The framework allows understanding the heterogeneity in the distributions, and the variation is attributed to spatial factors or other covariates. As well as introducing the models, it is discussed how identifiability, variable selection, and overfitting can be addressed within a Bayesian framework to provide interpretable inferential methods. Their use is illustrated in a range of applications.

Saturday 13.12.2025

10:05 - 12:10

Parallel Session B – CFE-CMStatistics 2025

CI015 Room BCB 307 DENSITY DATA ANALYSIS**Chair: Almond Stoecker****C0401: Functional principal component analysis for univariate and multivariate density data***Presenter:* **Karel Hron**, Palacky University, Czech Republic*Co-authors:* Adela Czolkova, Sonja Greven

The Bayes space provides a Hilbert space structure for the analysis of probability density functions (PDFs), equipping them with a geometry that respects their relative and constrained nature. A key tool in this framework is the centered logratio (clr) transformation, which establishes an isometric isomorphism between the Bayes space and the classical L^2 space. This enables the application of functional data analysis (FDA) techniques - in the context of dimension reduction, particularly functional principal component analysis (FPCA) - to both univariate and multivariate density data. In the multivariate setting, embedding PDFs in the Bayes space allows for an orthogonal decomposition into independent and interactive components; furthermore, the independent part can be decomposed into mutually orthogonal geometric marginals. This structure yields a deeper understanding of the sources of variation in multivariate densities and has direct consequences for interpreting the eigenfunctions and scores resulting from FPCA. It is shown that applying FPCA directly to multivariate densities is equivalent to applying multivariate FPCA to their decomposed form and that the resulting eigenfunctions and scores decompose accordingly. The theoretical results are illustrated with an application to empirical data, demonstrating the interpretability and practical value of this approach.

C0178: Density data analysis for densities observed via samples*Presenter:* **Sonja Greven**, Humboldt University of Berlin, Germany

In density data analysis, densities are in many settings the objects of interest and analysis, but are latent and only observed via samples. Common two-step approaches then first reconstruct densities using methods such as kernel density estimation or (compositional) splines, and ignore estimation uncertainty in the subsequent density data analysis. However, these approaches can be inaccurate, particularly if small or heterogeneous numbers of samples per density are available. The aim is to propose modeling individual draws from latent densities directly to incorporate all sources of uncertainty. This approach is illustrated for the cases of density principal component analysis as well as regression with densities as outcomes or covariates. To account for the constrained nature of densities, we base our approaches on Bayes spaces, which extend the Aitchison geometry for compositional data (discrete densities) to more general densities. Estimation can be based on (penalized) maximum likelihood estimation, in some cases requiring a (Monte Carlo) expectation maximization algorithm to handle latent densities. All approaches are illustrated with applications ranging from gender economics to climate research.

C0977: Nonparametric distributional inference using likelihoods as distributional data*Presenter:* **Karl Gerald van den Boogaart**, Helmholtz Zentrum Dresden Rossendorf e.V., Germany

One kind of distributional or density data are likelihoods. This contribution is concerned with a non-parametric Bayes estimation of a (multivariate) distribution from likelihood data through Bayes space methods. For likelihoods generated by complete point observations, the subspace of the Bayes space required for the posterior distribution of the underlying distribution is approximated by a relatively small subspace, only one dimension larger than the Bayes space of the original distribution. The situation, however, gets much more complicated if the data can no longer be represented by likelihoods of point observations, but starts to incorporate partially missing data or measurement errors. Such data can still be represented in terms of likelihoods and their Bayes space representation, but is no longer limited to this subspace. It is, however, possible to project those likelihoods onto the original subspace of the point observation likelihoods and still provide a reasonable Bayes inference that can be stored in a similar number of coefficients as the underlying multivariate distribution and provides relevant knowledge about the most relevant projection of the posterior distribution of the unknown distribution.

CO154 Room BCB G07 HiTEC: MACHINE LEARNING IN ECONOMETRICS**Chair: Juan Manuel Rodriguez-Poo****C1052: Structural inequalities in health across sub-Saharan Africa: A machine learning exploration of underlying determinants***Presenter:* **Mercedes Tejeria-Martinez**, University of Cantabria, Spain*Co-authors:* Vanesa Jorda, Jose Maria Sarabia

Health inequalities across sub-Saharan Africa (SSA), where disparities are exacerbated by limited access to healthcare and widespread socioeconomic inequities, represent critical barriers to equitable development. Advanced machine learning methodologies are leveraged to analyse health disparities in SSA countries using body mass index as a primary health indicator. Using demographic and health survey data, multiple machine learning algorithms are implemented to classify populations based on socioeconomic status, demographic profiles, and behavioural factors. The ensemble approach incorporates cross-validation model comparison techniques to optimise predictive accuracy and minimise algorithmic bias. The machine learning framework is designed to capture complex non-linear associations and interaction effects that conventional epidemiological approaches may miss. This computational approach aims to advance understanding of health inequality mechanisms in SSA, demonstrating how machine learning can enhance traditional public health research methodologies. The contribution to emerging digital health equity research is by establishing machine learning as a powerful tool for informing targeted, evidence-based interventions addressing systematic health disparities in developing regions.

C1053: A Gaussian process approach for testing varying coefficient models*Presenter:* **Luis Antonio Arteaga Molina**, Universidad de Cantabria, Spain*Co-authors:* Juan Manuel Rodriguez-Poo

The aim is to propose a Gaussian process (GP) approach for constancy testing in varying coefficient models. This methodology leads to a general unified framework of kernel-based tests having the following properties: (i) bootstrap tests are easy to implement in the presence of nuisance parameters (they are simple quadratic forms, and there is no need to re-estimate the nuisance parameters in each bootstrap replication); and (ii) the tests are valid under general conditions, including regularized estimators (e.g. Lasso) or parameters at the boundary of the parameter space. Neyman-orthogonal kernels are used, and a new asymptotic theory and a detailed local power analysis are developed. Monte Carlo experiments and a real data application illustrate the sensitivity of tests to the dimension of covariates and to the mean and covariance kernel of the GP.

C1001: Automatic debiased estimation with machine learning-generated regressors*Presenter:* **Telmo Perez**, University of the Basque Country (UPV/EHU), Spain*Co-authors:* Juan Carlos Escanciano

Many parameters of interest in economics and other social sciences depend on generated regressors. Examples in economics include structural parameters in models with endogenous variables estimated by control functions and in models with sample selection, treatment effect estimation with propensity score matching, and marginal treatment effects. More recently, machine learning (ML) generated regressors are becoming ubiquitous for these and other applications such as imputation with missing regressors, dimension reduction, including autoencoders, learned proxies, confounders and treatments, and for feature engineering with unstructured data, among others. The first general method is provided for valid inference with regressors generated from ML. Inference with generated regressors is complicated by the very complex expression for influence functions and asymptotic variances. Additionally, ML-generated regressors may lead to large biases in downstream inferences. To address these problems, automatic locally robust/debiased GMM estimators are proposed in a general three-step setting with ML-generated regressors. The

results are illustrated with treatment effects and counterfactual parameters in the partially linear and nonparametric models with ML-generated regressors. Sufficient conditions are provided for the asymptotic normality of the debiased GMM estimators, and their finite-sample performance is investigated through Monte Carlo simulations.

C1219: **Functional coefficient panel data models with endogenous selection: A flexible estimation approach**

Presenter: **Alexandra Soberon**, Universidad de Cantabria, Spain

Co-authors: Daniel Henderson, Juan Manuel Rodriguez-Poo, Taining Wang

The purpose is to introduce an innovative estimation approach for functional coefficient panel data models that simultaneously addresses sample selection and fixed effects challenges. Nearest neighbor differencing of smoothing variables is exploited to eliminate fixed effects without restrictive identification assumptions. The proposed two-stage technique seamlessly integrates nonparametric selection modeling with flexible functional coefficient estimation. The first stage captures complex selection mechanisms using advanced nonparametric techniques, while the second stage employs a sophisticated generalized local weighting scheme that estimates primary relationships while purging selection bias. The approach operates under weak regularity conditions while achieving superior computational efficiency compared to existing methods. Asymptotic properties are established, and it is demonstrated through Monte Carlo simulations that the estimators consistently outperform alternatives, delivering substantial improvements in bias reduction and precision.

C1525: **Financial attention and disclosure tone: Mixed frequency analysis**

Presenter: **Alev Atak**, METU, Turkey

The aim is to examine how demand-side financial attention and supply-side disclosure tone jointly relate to market volatility in a retail-dominated emerging market. A monthly Financial Attention Index (FAI) is constructed from Turkish-language Wikipedia pageviews via PCA; disclosure tone is measured by a FinBERT-based, probability- and market-cap-weighted NetTone Polarity Index (PI) from Borsa Istanbul annual reports (2016–2024). Mixed Data Sampling (MIDAS) regressions integrate indicators at native frequencies. Granger tests show that lagged FAI predicts volatility (not vice versa), and an event-driven difference-in-differences (DiD) around the February 2023 earthquake indicates a temporary strengthening of the attention–volatility link. By contrast, the tone–volatility association is negative in baseline models, whereas it becomes statistically insignificant after controlling for year fixed effects, consistent with tone primarily reflecting annual macroeconomic states rather than an independent within-year behavioral effect. Two temporal-direction tests – Granger causality and an event-driven DiD design – address concerns about simultaneity and clarify temporal precedence. These findings support the Wikipedia-based attention indices as scalable, real-time surveillance tools in retail-heavy emerging markets, where the FAI captures high-frequency behavioral dynamics, while disclosure tone reflects slower-moving macroeconomic conditions.

CO137 Room BCB G08 MACROECONOMIC RISK

Chair: Domenico Giannone

C0355: **Flexible priors and restrictions for structural vector autoregressions**

Presenter: **Francesca Loria**, Federal Reserve Board, United States

Co-authors: Christiane Baumeister, Junior Maih

The purpose is to introduce an innovative approach for estimating Bayesian vector autoregressions (VAR) in structural form, enhancing flexibility in incorporating various priors and identification strategies. The method accommodates zero, sign, and narrative restrictions, as well as identification via proxy variables, offering a unified framework for replicating prominent strategies in the VAR literature. Unlike existing methods, it directly elicits informative priors and restrictions on structural parameters, ensuring transparency and avoiding unintended beliefs. The approach is versatile, scalable to larger models, and eliminates the need for separate algorithms for different identification schemes. Additionally, the methodology allows for imposing (in-)equality restrictions on VAR parameters, providing a robust means to incorporate strong beliefs. Overall, this user-friendly framework addresses key challenges in the current literature, offering a valuable tool for empirical researchers, with the method accessible through the RISE toolbox.

C0365: **Scenario synthesis and macroeconomic risk**

Presenter: **Matteo Luciani**, Federal Reserve Board, United States

Co-authors: Domenico Giannone, Mike West, Tobias Adrian

Methodology is introduced to bridge scenario analysis and model-based risk forecasting, leveraging their respective strengths in policy settings. The Bayesian framework addresses the fundamental challenge of reconciling judgmental narrative approaches with statistical forecasting. Analysis evaluates explicit measures of concordance of scenarios with a reference forecasting model, delivers Bayesian predictive synthesis of the scenarios to best match that reference, and addresses scenario set incompleteness. This underlies systematic evaluation and integration of risks from different scenarios, and quantifies relative support for scenarios modulo the defined reference forecasts. The framework offers advances in forecasting in policy institutions that support clear and rigorous communication of evolving risks. Broader questions of integrating judgmental information with statistical model-based forecasts in the face of unexpected circumstances are also discussed.

C0786: **Forecasting macroeconomic risks in the UK**

Presenter: **Simon Lloyd**, Bank of England, United Kingdom

Co-authors: David Aikman, David Aikman, Rhys Bidder, Giulia Mantoan, Simone Maso, Aditya Mori, Matthew Tong

Statistical metrics of UK macroeconomic risks are constructed, building quantile-regression density forecasts for inflation and GDP growth. The models account for the UK's position as a small-open economy and capture time variation in tail risks, owing to variation in economic and financial conditions, in addition to changes in central forecasts. To highlight how the statistical models of macroeconomic risk can provide valuable real-time signals about the balance of risks, a battery of tests is used to compare the predictive distributions to those captured in the Bank of England's fan charts. The fitted densities for growth are at least as well calibrated as the fan charts, outperforming the fans in terms of relative accuracy. Although the estimates for inflation perform similarly to the fans over the last two decades overall, they do better capture economic narratives when inflation deviates from the target. These tools can contribute to a broader suite for quantifying macroeconomic risks in the UK, with regular evaluation of density forecasts necessary to ensure that the toolkit remains fit for purpose as the constellation of shocks hitting the UK economy evolves.

CO094 Room BCB G09 ADVANCES IN MACRO- AND FINANCIAL ECONOMETRICS

Chair: Toshiaki Watanabe

C0183: **Time-varying local projections with application to trade volume analysis**

Presenter: **Jouchi Nakajima**, Hitotsubashi University, Japan

A general approach to dynamic modeling is discussed using the local projection (LP) method. Previous studies have proposed time-varying (TV) parameters in LPs; However, they did not address possible variations in error variances. Overlooking this could introduce significant bias in the estimate of the TV parameter, and consequently, the estimated impulse response. An estimation strategy is developed for LPs with stochastic volatility (SV), and the importance of SV inclusion is illustrated using simulated data. An application to a topical macroeconomic time-series analysis illustrates the benefits of the proposed approach in terms of improved predictions.

C0211: **Multi-view dynamic network modeling**

Presenter: **Mike So**, The Hong Kong University of Science and Technology, Hong Kong

Co-authors: Shun Hin Chan, Amanda Chu

A flexible multi-view dynamic network model is developed using a regression-like structure, incorporating exogenous and endogenous variables from the lagged networks to model edge changes. The model does not rely on latent space, simplifying network estimation and prediction. Furthermore, it integrates a multi-view feature to represent various relationship types at each time point. The proposed model offers an intuitive interpretation of the estimation. Bayesian model averaging method is also applied to predict networks.

C0259: The determinants of liquidity in the Japanese government bond markets: An interpretable machine learning approach

Presenter: **Satoko Kojima**, The Bank of Japan, Japan

Co-authors: Toshiyuki Sakiyama

Liquidity in government bond markets is critical for the functioning of financial markets. The aim is to study what and how bond features drive the liquidity measured by price dispersion by applying machine learning approaches to high-granularity data from the Bank of Japan financial network system. The main findings are threefold. First, the decomposition of the liquidity indicator into bond features reveals that the historical volatility of benchmark prices has been the main driver of the liquidity indicator, while the contributions of the share of non-clearing participants' transactions and the share of the central bank's transactions and holdings have increased since around 2021. Second, some bond features affect the liquidity indicator non-linearly. For bond features such as the share of foreign financial institutions' transactions, the number of trading financial institutions, and the share of the central bank's holdings, the liquidity indicator improves as the values of these bond features increase, but it deteriorates once they exceed certain thresholds. Third, bond features such as maturity and the number of trading financial institutions interact strongly with other bond features in a way that changes how the interacted bond features affect the liquidity indicator.

C0487: Bayesian forecasting of tail risk in cryptocurrency markets

Presenter: **Cathy W-S Chen**, Feng Chia University, Taiwan

Co-authors: Ying-Lin Hsu, Po-Hui Chen

Cryptocurrencies exhibit high volatility, emphasizing the importance of accurately measuring tail risk in their markets. A threshold-switching mechanism is incorporated into Taylor's ES-CAViaR models that unveil features such as asymmetry and jump phenomena. These enhancements effectively capture the diverse tail risks of cryptocurrencies while enabling the simultaneous forecasting of both value-at-risk (VaR) and expected shortfall (ES). The proposed models incorporate two types of functions to address the VaR and ES nexus, with the option to use the rolling standard deviation of returns as a short-term volatility proxy as a regressor. The parameters and forecast tail risk are estimated within a Bayesian framework. Taking the two largest cryptocurrencies by market capitalization, Bitcoin and Ethereum, the one-step-ahead forecasting performance is assessed over a four-year out-of-sample period using a rolling window approach. The comparative results from backtests and five scoring functions among eight competing models support the conclusion that models with a threshold mechanism capture the tail risk of cryptocurrencies more accurately than other risk models.

C0615: Loss-based Bayesian sequential prediction of value-at-risk with a long-memory and non-linear realized volatility model

Presenter: **Richard Gerlach**, University of Sydney, Australia

Co-authors: Chao Wang, Minh-Ngoc Tran, Rangika Peiris

A long-memory and non-linear realized volatility model class is proposed for direct value-at-risk (VaR) forecasting. This model, referred to as RNN-HAR, extends the heterogeneous autoregressive (HAR) model, a framework known for efficiently capturing long memory in realized measures, by integrating a recurrent neural network (RNN) to handle the non-linear dynamics. Quantile loss-based generalized Bayesian method with sequential Monte Carlo is employed for model estimation and sequential prediction in RNN-HAR. The empirical analysis is conducted using daily closing prices and realized measures with around 12 years of data till 2022, covering 31 market indices. The proposed model's one-step-ahead VaR forecasting performance is compared against a basic HAR model and its extensions. The results demonstrate that the proposed RNN-HAR model consistently outperforms all other models considered in the study.

| | | |
|-------------------------------|---|-----------------------|
| C0296 Room Virtual R01 | ADVANCES IN PANEL DATA AND TIME SERIES MODELS FOR FINANCIAL ECONOMETRICS | Chair: Fei Liu |
|-------------------------------|---|-----------------------|

C0309: Time-varying forecasting regressions augmented by a dynamic nonstationary factor structure

Presenter: **Tingting Cheng**, Nankai University, China

The purpose is to introduce a time-varying predictive regression augmented by nonstationary common factors, considering both static and dynamic factor models. An easy-to-implement sieve-based estimation method is proposed to estimate the unknown time-varying coefficients in the predictive regression. Large sample theories are established for those estimators and also examine their finite-sample performance by simulation studies. The advantages of this new model are further demonstrated through an empirical application to U.S. inflation forecasting.

C0260: Augmented LASSO with textual analysis in the financial market

Presenter: **Shuyi Ge**, University of Nankai, China

An augmented lasso (ALasso) framework that integrates textual information is proposed to enhance asset return prediction. Specifically, firm linkages and word-level sentiment are extracted from news reports, and these signals are embedded into the Lasso estimation process. Firm linkage strength is used to impose discriminative penalties across firms, enhancing cross-firm predictability, while word sentiment serves as a directional prior to constrain the signs of textual predictors. The oracle properties of ALasso are theoretically established under varying levels of news information quality. Empirically, ALasso demonstrates superior performance in predicting returns in the A-share market.

C0555: Robust estimation and inference for high-dimensional panel data models

Presenter: **Yayi Yan**, Shanghai University of Finance and Economics, China

Co-authors: Jiti Gao, Fei Liu, Bin Peng

The relevant literature is provided with a complete toolkit for conducting robust estimation and inference about the parameters of interest involved in a high-dimensional panel data framework. Specifically, (1) non-Gaussian, serially and cross-sectionally correlated and heteroskedastic error processes are allowed for, (2) an estimation method is developed for a high-dimensional long-run covariance matrix using a thresholded estimator, and (3) the number of regressors is also allowed to grow faster than the sample size. Methodologically and technically, two Nagaev types of concentration inequalities are developed: One for a partial sum and the other for a quadratic form, subject to a set of easily verifiable conditions. Leveraging these two inequalities, a non-asymptotic bound is derived for the LASSO estimator, achieving asymptotic normality via the node-wise LASSO regression, and establishing a sharp convergence rate for the thresholded heteroskedasticity and autocorrelation consistent (HAC) estimator. The practical relevance of these theoretical results is demonstrated by investigating a high-dimensional panel data model with interactive effects. Moreover, extensive numerical studies are conducted using simulated and real data examples.

C0191: Panel data estimation and inference: Homogeneity versus heterogeneity

Presenter: **Fei Liu**, Nankai University, China

An underlying data-generating process that allows for different magnitudes of cross-sectional dependence is defined, along with time series autocorrelation. This is achieved via high-dimensional moving average processes of infinite order (HDMA(∞)). The setup and investigation integrate and enhance homogeneous and heterogeneous panel data estimation and testing in a unified way. To study HDMA(∞), the Beveridge-Nelson decomposition is extended to a high-dimensional time series setting, and a complete toolkit set is derived. Homogeneity versus heterogeneity is examined using Gaussian approximation, a prevalent technique for establishing uniform inference. For post-testing inference, central limit theo-

rems are derived through Edgeworth expansions for both homogeneous and heterogeneous settings. Additionally, the practical relevance of the established asymptotic properties is showcased by revisiting the common correlated effects (CCE) estimators and a classic nonstationary panel data process. Finally, the theoretical findings are verified via extensive numerical studies using both simulated and real datasets.

CO036 Room BCB 206 CFE SESSION: A TRIBUTE TO H. PESARAN I
Chair: Natalia Bailey
C0424: Analysis of multiple long run relations in panel data models with applications to financial ratios

Presenter: **Alexander Chudik**, Federal Reserve Bank of Dallas, United States

Co-authors: Hashem Pesaran, Ron Smith

A new methodology is provided for the analysis of multiple long-run relations in panel data models where the cross-section dimension, n , is large relative to the time series dimension, T . For panel data models with large n , researchers have focused on panels with a single long-run relationship. The main difficulty has been to eliminate short-run dynamics without generating significant uncertainty for the identification of the long run. This problem is overcome by using non-overlapping sub-sample time averages as deviations from their full-sample counterpart and estimating the number of long-run relations and their coefficients using eigenvalues and eigenvectors of the pooled covariance matrix of these sub-sample deviations. This procedure is referred to as pooled minimum eigenvalue (PME), and it is shown that it applies to unbalanced panels generated from general linear processes with interactive stationary time effects and does not require knowing long-run causal linkages. To best of knowledge, no other estimation procedure exists for this setting. The PME estimator is shown to be consistent and asymptotically normal as n and $T \rightarrow \infty$ jointly, such that $T \approx n^d$, with $d > 0$ for consistency and $d > 1/2$ for asymptotic normality. Extensive Monte Carlo studies show that the number of long-run relations can be estimated with high precision, and the PME estimates of the long-run coefficients show small bias and RMSE and have good size and power properties.

C0498: Forecasting with heterogeneous spatiotemporal models

Presenter: **Cynthia Yang**, Florida State University, United States

The purpose is to investigate the forecasting performance of dynamic spatial panel data models that incorporate heterogeneous coefficients and latent common factors. A comprehensive comparison of the predictive accuracy of univariate versus multivariate models is conducted, with and without accounting for spatial interdependence, across various forecasting horizons. Monte Carlo simulations are employed to assess the finite sample properties of competing models and estimators. An empirical application focused on house price forecasting demonstrates the practical advantages of the spatiotemporal models.

C0594: Stochastic search selection for heterogeneous panel data models

Presenter: **Andreas Pick**, Erasmus University Rotterdam, Netherlands

The aim is to present a method for selecting variables and determining parameter heterogeneity in Bayesian hierarchical panel data models. Mixture distributions are used as priors for the mean and the variance of the individuals' parameters. Selection indicators determine the best-fitting component of each mixture distribution and indicate whether the mean parameter is non-zero and whether the parameters are heterogeneous. The method is applied to panel data sets on inflation of US CPI sub-indices and house price inflation across US metropolitan statistical areas.

C0245: Continuously updating GMM estimator: Bridging theory and practice

Presenter: **Mahrad Sharifvaghefi**, University of Pittsburgh, United States

Co-authors: Whitney Newey, Marcelo J Moreira

The continuously updating generalized method of moments (CU-GMM) estimator has appealing theoretical properties. However, it lacks a closed-form expression for HAC errors and requires solving a non-convex optimization problem. It is demonstrated that the objective function of the CU-GMM estimator is a ratio of polynomials. Consequently, the optimization problem is translated into finding roots for polynomial equations. The analysis reveals the sensitivity of the CU-GMM estimator to different optimization methods, showcasing significantly reduced dispersion compared to prior literature.

C0548: Robust tests for quantile-based directional predictability

Presenter: **Natalia Bailey**, Monash University, Australia

Co-authors: Bonsoo Koo, Myung Hwan Seo

Quantile dependence analysis within and between time series offers valuable insights into the dynamics of economic and financial data, particularly in capturing tail dependencies that traditional correlation measures often overlook. The quantilogram and cross-quantilogram are statistical tools designed to capture dependencies across different segments of the distribution of a time series and the joint distribution of two time series, respectively. This makes them particularly useful tools in financial econometrics, especially when analyzing extreme events. Test statistics of no quantile directional predictability are developed, which are also robust to complex volatility structures and non-linearities commonly characterizing financial time series. The asymptotic properties of the self-normalized statistics are established, and it is demonstrated that they yield valid confidence bands which produce correct size and satisfactory power without requiring the estimation of long-run variance or choice of tuning parameters. Monte Carlo simulation results are encouraging. The empirical application evaluates the evolution of build-up in systemic risk in the US financial market over the period 1993-2024 as measured by 5-year rolling estimates of the (cross-)quantilograms of risk-adjusted returns for three selected securities (JPM, MS, AIG). A significant increase in directional predictability at the fifth quantile is uniformly detected during the Global Financial Crisis.

CO219 Room BCB 207 STATISTICAL METHODS FOR INSURANCE
Chair: Gabriella Piscopo
C0536: Shapley risk sharing in peer-to-peer insurance

Presenter: **Susanna Levantesi**, Sapienza University of Rome, Italy

Co-authors: Gian Paolo Clemente, Gabriella Piscopo

Peer-to-peer (P2P) insurance is an innovative model that leverages digital technology to connect individuals with similar insurance needs, forming a pool to share risks. A P2P insurance framework is introduced in which participants pay an ex-ante contribution determined by the Shapley value, where the value-at-risk of the difference between the expected and realized total loss of the network (i.e., profit and loss) is used as the risk measure. Under standard assumptions commonly used in non-life insurance for the aggregate claim amount, closed-form expressions are derived for the Shapley value. The model includes a cashback mechanism to ensure that all participants contribute equally to covering realized losses. The practical implementation of the model is demonstrated by applying it to a portfolio of motor other damage insurance policies.

C0546: A clustering approach to assess house price risk in reverse mortgages: Insights on the Italian market

Presenter: **Giulia Magni**, Sapienza University of Rome, Italy

Co-authors: Emilia di Lorenzo, Alba Roviello, Marilena Sibillo

Although positive progress, aging of the population and increasing life expectancy have intensified poverty and inequality, even in advanced economies. The reverse mortgage is a supplementary pension product for elderly homeowners with valuable property but limited liquidity ('house rich and cash poor'). It allows them to borrow money using their home as a guarantee, maintaining the legal rights to the property. Debt repayment is deferred to heirs, who can settle the loan and reclaim the property. Heirs are ensured by the non-negative equity guarantee, which limits repayment to the actual value of the property. Unlike other countries, reverse mortgage is not widely developed in Italy, although they fit well into the context of aging and a high concentration of real estate ownership. The risks related to mortality, interest rates, and the volatility of

house prices over time hinder its issuance. Above all, the house price risk undertakes a predominant role. Territorial analysis to differentiate areas with similar characteristics enables better management of this criticality. Exploiting geographically granular Italian real estate market values, a clustering procedure is adopted to account for territorial characteristics in the contractual specifications. This approach allows for performing risk segmentation and accurately representing the actual risk associated with each territorial zone.

C0733: The ESG score and the sustainable development: A machine learning analysis

Presenter: **Gabriella Piscopo**, University of Naples Federico II, Italy

Co-authors: Susanna Levantesi, Kevyn Stefanelli

In line with the values of a sustainable economy, companies are progressively implementing strategies that aim to balance profitability with environmental, social, and governance (ESG) commitments. The financial sector's heightened sensitivity to climate and environmental risks accentuates the imperative of advancing sustainable investment practices. Within this framework, sustainability, integrating ESG factors, stands as a central strategic focus. The aim is to investigate the relation between ESG score and sustainable practice of some listed companies. To this end, the results of the classical regression framework are compared with those of advanced machine learning techniques, including random forest and gradient boosting machine algorithms.

C0737: Data-driven policy design for parametric insurance on the energy production of a photovoltaic system

Presenter: **Fabio Baione**, Sapienza University of Rome, Italy

Co-authors: Emiliano Valente

Climate change represents one of the most pressing and complex challenges of the time. The objective is to propose the formulation of a parametric insurance policy to cover the risks arising from low electricity production from a photovoltaic system. Such policies can serve as an incentive to adopt renewable energy sources, encouraging their use in contexts such as powering private homes or small towns. To this end, the global horizontal irradiance (GHI) variable is used as a reference index, which represents the solar irradiance absorbed by a surface parallel to the Earth's surface. Assuming a Kumaraswamy distribution for solar irradiance in each hour, each hourly variable was interpreted as a marginal distribution belonging to a multivariate distribution function representing the entire day, modeled using a Gaussian copula.

C0909: Modeling weekly temperature-mortality dynamics in Italy (2011/2024) with EVT-based analysis

Presenter: **Francesca Serra**, La Sapienza, Italy

The aim is to investigate the relationship between extreme temperatures and weekly elderly mortality in Italy from 2011 to 2024 using a two-step approach. First, regional indicators of extreme heat (T90) and cold (T10) are developed by comparing weekly temperature values with historical thresholds from a 1985-2000 reference period. Weekly mortality rates are computed for five elderly age groups using official data. Initial findings suggest non-linear associations: Cold extremes (T10) are more strongly linked to excess mortality in northern and central regions, while heat extremes (T90) are more relevant in the South. To remove trends and seasonality, ARIMA and ARIMA-GARCH models are applied to both temperature and mortality series, and standardized residuals are extracted. In the second step, extreme value theory is used to analyze the joint tail behavior of residual temperature and mortality. A Peaks-over-Threshold method combined with bivariate copulas models the dependence in the extremes. Results reveal more frequent and intense co-extremes between high temperature anomalies and mortality, especially among the oldest groups in southern regions. Cold-related co-extremes appear weaker and more regionally dispersed. These findings highlight the value of modeling joint extremes to better understand climate-related health risks and inform adaptation strategies for vulnerable populations.

CO266 Room BCB 208 MODERN STATISTICAL METHODS FOR THE ANALYSIS OF COMPLEX REAL DATA

Chair: Joanna Janczura

C0491: Distinguishing anomalous diffusion: A statistical approach to parameter characterization

Presenter: **Agnieszka Wylomanska**, Wroclaw University of Science and Technology, Poland

Co-authors: Katarzyna Maraj-Zygmunt, Aleksandra Grzesiek, Diego Krapf

Anomalous diffusion describes processes where a particle's mean squared displacement scales non-linearly with time. This behavior, common in complex systems (like biological cells), goes beyond standard diffusion. Traditional models like fractional Brownian motion (FBM) and scaled Brownian motion (SBM) assume a constant anomalous diffusion exponent, which limits their ability to capture dynamics with varying anomalous parameters. To overcome this, FBM with random exponents (FBMRE) and SBM with random exponents (SBMRE) were developed. This research introduces a universal statistical testing framework to differentiate between anomalous diffusion models having constant versus random anomalous exponents. It uses time-averaged statistics and their ratios. This methodology broadly applies to constant vs. random anomalous diffusion scenarios, with its effectiveness depending on chosen statistics, time lags, and process properties. This is demonstrated through simulations (using a two-point exponent distribution) and real-world data analysis.

C1048: Multidimensional fractional Brownian motion: From construction to analysis

Presenter: **Michał Balcerek**, Wroclaw University of Science and Technology, Poland

The aim is to introduce a novel construction of two-dimensional fractional Brownian motion (2d fBm) that explicitly captures both anisotropic scaling and cross-component dependencies using a matrix-valued Hurst operator and correlated Gaussian noise. Two distinct models are proposed: Causal and well-balanced 2d fBm, differing in their cross-covariance structures. The theoretical properties of such models (both in time and frequency domains) are introduced, and those are compared to the ones obtained from the numerical simulations. Particular attention is given to how variations in the Hurst parameters and noise correlation structures manifest in observed behavior, and how those influence data analysis.

C0507: MIDAST as a novel approach to multivariate data segmentation - smart brain monitoring

Presenter: **Justyna Witulska**, Wroclaw University of Science and Technology, Poland

Co-authors: Marta Hendler, Magdalena Kasproicz, Marek Czosnyka, Ireneusz Jablonski, Agnieszka Wylomanska

The research addresses the problem of identifying distributional changes in multivariate non-Gaussian data. A novel and general methodology, called MIDAST, is introduced for fusion-based segmentation of multivariate data using statistical tests (e.g., the Kolmogorov-Smirnov test, a maximum mean discrepancy-based test, and a kernel-based test). The proposed approach is evaluated against baseline methods, including E-Divisive and Kernel Change-point Analysis (KCPA), with segmentation accuracy and computational complexity as key performance indicators. The methodology is tested through computer simulations, using multivariate sub-Gaussian and multivariate Student's *t* distributions, under varying degrees of correlation, degrees of freedom, and stability indices. MIDAST, enhanced by a windowing mechanism, enables the detection of one or multiple change points that indicate shifts in the statistical properties of the system. A real-world application demonstrates the ability of the method to reduce invasiveness in intracranial hypertension event detection by identifying structural changes in multivariate temporal patterns. The proposed approach facilitates more accurate and adaptive monitoring of complex systems governed by evolving statistical dynamics.

C1192: Statistical inference with the modified Greenwood statistic for univariate and multivariate heavy-tailed models

Presenter: **Marek Arendarczyk**, University of Wroclaw, Poland

Co-authors: Marek Arendarczyk, Tomasz Kozubowski, Anna Panorska, Katarzyna Skowronek, Agnieszka Wylomanska

The Greenwood statistic and its modifications play an important role in modern statistical methodology for extremes, heavy-tailed distributions, and non-Gaussian models. Originally introduced for testing exponentiality, the statistic has since been developed into a versatile tool with applications in extreme value analysis and the study of clustering and heterogeneity. Recent results demonstrate stochastic ordering of the Greenwood statistic with respect to the tail index, which makes it suitable for building tests and confidence intervals in families of distributions relevant for extreme value

theory. Further developments adapt the modified Greenwood statistic to symmetric alpha-stable and Student's t families of distributions. In addition, a multivariate generalization has been proposed, extending the use of the Greenwood framework to vector-valued samples, including the sub-Gaussian case, multivariate Pareto, and multivariate Student's t families of distributions. Consequently, the modified Greenwood statistic enables effective testing of multivariate Gaussianity, identification of infinite variance, and discrimination between heavy- and light-tailed multivariate models. The mathematical framework is discussed in stochastic ordering results, simulation studies, and real data applications, emphasizing the modified Greenwood statistic as an important and efficient tool in statistical inference.

C0918: Inference from joint distributions: Product of random variables with an application to energy market

Presenter: **Joanna Janczura**, Wrocław University of Science and Technology, Poland

Co-authors: Agnieszka Wylomanska, Andrzej Puc

Multivariate analysis is a cornerstone of modern science, enabling comprehensive investigations of complex phenomena, where multiple variables interact. Despite the critical role of joint distributions of variables, in certain scientific applications, it's essential to study the product of variables rather than their joint distribution. Among many others, these are economic variables with random discount factors or tax rates, transaction values, costs of prediction errors, or variables with a random scale parameter. The aim is to study the distributional properties of a product of random variables as well as the product of vector autoregression model components. For the introduced time series, general formulas are derived for the autocovariance function, and its properties are studied for different cases of cross-dependence structure. The theoretical results are then illustrated using simulations and are applied to an electricity market case study in which the financial cost of balancing load prediction errors is analyzed after the day-ahead market settlement and prior to delivery. The considered approach yields a model that is consistent for multivariate time series as well as their product, and, on the other hand, can be an economically grounded alternative for statistical evaluation of the load (or price) forecast accuracy.

CO208 Room BCB 209 INFERENCE AND CHANGE DETECTION FOR DEPENDENT DATA

Chair: Ansgar Steland

C0563: Asymptotic properties of change point detection in high-dimensional data with a strongly spiked eigenvalue structure

Presenter: **Kento Egashira**, Tokyo University of Science, Japan

Co-authors: Kazuyoshi Yata, Makoto Aoshima

The aim is to consider detecting change points in high-dimensional data with limited sample sizes, particularly under a strongly spiked eigenvalue (SSE) model. A multivariate CUSUM-type statistic is introduced, designed to compare the means before and after each potential change point. Unlike many existing techniques, the approach avoids imposing sparsity constraints. The asymptotic behavior of the proposed statistic is investigated under the null hypothesis of no structural change, and the consistency of the corresponding change-point estimator is demonstrated under mild regularity conditions. In addition, the null distribution of the test statistic is derived specifically under the SSE setting. Extensive numerical experiments confirm the practical utility and robustness of the proposed methodology.

C0698: Weak convergence of the partial sum of $I(d)$ process to a fractional Brownian motion in finite interval representation

Presenter: **Junichi Hirukawa**, Nanzan University, Japan

Co-authors: Kou Fujimori

An integral transformation that changes a fractional Brownian motion to a process with independent increments has been given. A representation of a fractional Brownian motion through a standard Brownian motion on a finite interval has also been given. On the other hand, it is known that the partial sum of the discrete time fractionally integrated process ($I(d)$ process) weakly converges to a fractional Brownian motion in infinite interval representation. The weak convergence of the partial sum of the $I(d)$ process to a fractional Brownian motion is derived in the finite interval representation.

C0714: A Gaussian approximation result for weakly dependent random fields using dependency graphs

Presenter: **Dennis Loboda**, RWTH Aachen University, Institute of Statistics, Germany

Non-stationary random fields under the physical dependence measure are investigated. In particular, the objective is to study the maximum of local averages given an increasing bandwidth under expanding-domain asymptotics. By defining suitable vectors based on the studied random field, it becomes possible to use the concept of dependency graphs known from time series analysis. This leads to an approximation result for the maximum of local averages through a Gaussian random field, which preserves the covariance structure.

C0818: Sliced-Wasserstein distance based change detection with sequential empirical processes

Presenter: **Florian Scholze**, RWTH Aachen University/ University of Bamberg, Germany

Co-authors: Fabian Mies, Ansgar Steland

The purpose is to study the problem of detecting changes in the marginal distributions of a multivariate time series with a CUSUM-type detector statistic based on the (maximum-) sliced-Wasserstein distance. From a theoretical point of view, this projection-based approach has two appealing properties. Firstly, unlike the family of Wasserstein distances, it does not suffer from the curse of dimensionality, and secondly, by means of the Kantorovich duality, asymptotic properties of the detector statistic can be derived from results for function-indexed sequential empirical processes for nonstationary time series. A new (bootstrap-) functional central limit theorem is presented for sequential empirical processes and its application to the given testing problem. Practical implications, limitations, and possible extensions are discussed.

C0968: Quantitative inference about the variation of a function

Presenter: **Fabian Mies**, Delft University of Technology, Netherlands

Co-authors: Holger Dette

When performing nonparametric inference about a regression function, quantitative assessments are usually either point estimates, pointwise confidence intervals, or simultaneous confidence bands. While this gives some impression of the shape of the signal, it does not directly yield statistical statements about the variability of the function within the band. The aim is to quantify the variation of a function in terms of (a) its range, and (b) its total variation. Based on multiscale test statistics, simultaneous confidence intervals are constructed for these variation measures, as well as for the modulus of continuity of the regression function. For the special case of step functions, the methods yield novel multiscale tests for relevant change-points.

CO133 Room BCB 210 STATISTICS FOR STOCHASTIC PROCESSES AND THEIR APPLICATIONS

Chair: Markus Bibinger

C0774: Parametric estimation for weak fractional time series

Presenter: **Tetsuya Takabatake**, The University of Osaka, Japan

Parametric estimation is considered for fractional time series models with general memory parameters and an innovation process satisfying the weak white noise condition. Estimation is performed via a conditional-sum-of-squares estimator, based on truncating the infinite-order autoregressive representation of the observed fractional time series, covering both stationary and non-stationary ranges of the memory parameter. Consistency and asymptotic normality are established under minimal conditions on the innovation process. Moreover, time permitting, the estimator's performance is also demonstrated through simulations.

C1059: Prediction of linear fractional stable motions using codifference

Presenter: **Matthieu Garcin**, ESILV, France

Co-authors: Karl Sawaya, Thomas Valade

The linear fractional stable motion (LFSM) extends the fractional Brownian motion (fBm) by considering α -stable increments. A method is proposed to forecast future increments of the LFSM from past discrete-time observations, using the conditional expectation when $\alpha > 1$ or a semimetric projection otherwise. It relies on the codifference, which describes the serial dependence of the process, instead of the covariance. Indeed, covariance is commonly used for predicting an fBm, but it is infinite when $\alpha < 2$. Some theoretical properties of the method and of its accuracy are studied and both a simulation study and an application to real data confirm the relevance of the approach. The LFSM-based method outperforms the fBm, when forecasting high-frequency FX rates. It also shows a promising performance in the forecast of time series of volatilities, decomposing properly, in the fractal dynamic of rough volatilities, the contribution of the kurtosis of the increments, and the contribution of their serial dependence. Moreover, the analysis of hit ratios suggests that, besides independence, persistence, and antipersistence, a fourth regime of serial dependence exists for fractional processes, characterized by a selective memory controlled by a few large increments.

C0338: Estimating the Hurst parameter via ordinal pattern distributions

Presenter: Alexander Schnurr, University Siegen, Germany

The ordinal structure of long-range dependent time series is analyzed. To this end, so-called ordinal patterns are used, which describe the relative position of consecutive data points. Two estimators are provided for the probabilities of ordinal patterns and prove limit theorems in different settings, namely stationarity and (less restrictive) stationary increments. In the second setting, a Rosenblatt distribution in the limit is encountered. More general limit theorems are proven for functions with Hermite rank 1 and 2. The limit distribution is derived for an estimation of the Hurst parameter H if it is higher than $3/4$. Thus, the theorems complement results for lower values of H , which can be found in the literature.

C0402: The trade-off between model flexibility and accuracy of the expected threat model in football

Presenter: Jakob Soehl, Delft University of Technology, Netherlands

Co-authors: Koen van Arem, Mirjam Bruinsma, Geurt Jongbloed

With an average football (soccer) match recording over 3,000 on-ball events, effective use of this event data is essential for practitioners at football clubs to obtain meaningful insights. Models can extract more information from this data, and explainable methods can make them more accessible to practitioners. The expected threat model has been praised for its explainability and offers an accessible option. However, selecting the grid size is a challenging key design choice that has to be made when applying the expected threat model. Using a finer grid leads to a more flexible model that can better distinguish between different situations, but the accuracy of the estimates deteriorates with a more flexible model. Consequently, practitioners face challenges in balancing the trade-off between model flexibility and model accuracy. The expected threat model is analyzed from a theoretical perspective, and simulations are performed based on the Markov chain of the model to examine its behavior in practice. The theoretical results establish an upper bound on the error of the expected threat model for different flexibilities. Based on the simulations, a more accurate characterization of the model's error is provided, improving over the theoretical bound. Finally, these insights are converted into a practical rule of thumb to help practitioners choose the right balance between the model flexibility and the desired accuracy of the expected threat model.

C1112: Bayesian nonparametrics for semi-linear stochastic PDEs

Presenter: Randolph Altmeyer, Imperial College London, United Kingdom

Co-authors: Sascha Gaudlitz

The aim is to consider the nonparametric estimation of the reaction function in a semi-linear stochastic partial differential equation (SPDE) from observing a trajectory of the solution continuously over a finite time interval. Given a Gaussian process prior, Bayesian posterior contraction rates are derived in a novel asymptotic regime: The spatial domain grows while the time horizon remains fixed. In this setting, the solution of the SPDE converges to a stationary process and is spatially ergodic. This allows for proving concentration inequalities of functionals along spatial averages of the solution. The proofs rely on the Clark-Ocone formula from Malliavin Calculus and precise bounds on the marginal densities of the SPDE. The posterior contraction rates are minimax-optimal. A nonparametric Bernstein-von Mises theorem is further proven for the posterior distribution.

CO121 Room BCB 211 BAYESIAN ECONOMETRICS

Chair: Yasuhiro Omori

C0862: Dynamic Bayesian regression quantile synthesis for forecasting outlook-at-risk

Presenter: Genya Kobayashi, School of Commerce, Meiji University, Japan

Co-authors: Yuta Yamauchi, Shonosuke Sugawara, Dongu Han

The aim is to provide a Bayesian approach to accurate quantile forecasting for time series data through the Bayesian predictive synthesis. The proposed dynamic Bayesian regression quantile introduces predictions from the agent quantile predictive models as latent factors and lets the weights for the agent models vary across time, constituting a dynamic latent factor model for quantiles. It is also considered to extend the model for quantile prediction of multiple time series data by introducing an additional factor structure to the synthesis weights. The performance of the proposed approach is demonstrated using the US inflation rate and GDP growth rates for some developed countries.

C0932: Quantile forecasts with stochastic volatility models using realized quantile measures

Presenter: Yuta Yamauchi, Nagoya University, Japan

Co-authors: Genya Kobayashi, Yasuhiro Omori

The aim is to improve the forecasting accuracy of value-at-risk (VaR) by using information on quantiles derived from high-frequency data, based on stochastic volatility models. By incorporating realized volatility, the model enhances volatility forecasting while simultaneously utilizing auxiliary information from high-frequency data related to the quantiles. A dynamic model is constructed for both the asset return quantiles and volatility as latent processes, enabling high-frequency information to be effectively incorporated through the observation equations. An empirical analysis of asset return data is conducted to assess the predictive performance of the proposed model relative to existing realized stochastic volatility models.

C0997: Flexible skewness modeling in stochastic volatility: Generalized Fernandez-Steel distribution approach

Presenter: Tomoya Yano, The University of Tokyo, Graduate School of Economics, Japan

Co-authors: Yasuhiro Omori

The focus is on a stochastic volatility model with error terms following a generalized Fernandez-Steel distribution. Financial return data often exhibit characteristics such as heavy tails and skewness, which can be partially captured by the standard Fernandez-Steel distribution. However, by extending the distribution to allow connections at arbitrary quantile points, not just at the origin, the model gains greater flexibility. This enhancement is expected to improve the accuracy of forecasts for risk measures such as value-at-risk (VaR) and expected shortfall (ES).

C0998: Dynamic factor stochastic volatility in mean model

Presenter: Daichi Hiraki, University of Tokyo, Japan

Co-authors: Yasuhiro Omori

A stochastic volatility in mean (SVM) model is developed within a dynamic factor framework to capture common movements in macroeconomic variables under time-varying uncertainty. Motivated by theoretical considerations in macro-finance, the model allows conditional volatility to directly affect the conditional mean through a volatility-in-mean component. This feature enables the model to account for time-varying risk premiums that are otherwise difficult to capture in standard factor stochastic volatility models. The model is estimated using Bayesian Markov chain Monte Carlo methods and applied to quarterly U.S. macroeconomic data from the FRED-QD dataset. The empirical results illustrate how

the SVM structure can be embedded in a latent factor setting to study macroeconomic dynamics under uncertainty, providing a basis for future forecasting and structural analysis.

C1100: **Two-sample comparison through additive tree models for density ratios**

Presenter: **Naoki Awaya**, Waseda University, Japan

Co-authors: Li Ma, Yuliang Xu

The density ratio is an effective summary of the difference between two distributions. The aim is to propose additive tree ensembles for the density ratio, along with efficient algorithms for training these models based on i.i.d. samples from the distributions. A loss function is introduced called the balancing loss under which such models can be trained from both an optimization perspective that parallels tree boosting and from a (generalized) Bayesian perspective that parallels Bayesian additive regression trees (BART). For the former, two boosting algorithms are presented: One based on forward-stage-wise fitting and the other on gradient boosting for computing a single estimate for the density ratio function. For the latter, it is shown that due to its resemblance to an exponential family kernel, the new loss can serve as a pseudo-likelihood for which conjugate priors exist, thereby enabling effective generalized Bayesian inference on the density ratio using the backfitting sampler for BART. This allows generalized Bayesian uncertainty quantification on the inferred density ratio, which is critical but often unaddressed in modern applications involving two-sample comparison. The application of the method is demonstrated in a case study involving assessing the quality of generative models for microbiome compositional data.

CO180 Room BCB 212 COMPUTATIONAL METHODS IN STATISTICS

Chair: Martina Amongero

C0439: **Zero-order parallel sampling**

Presenter: **Francesco Pozza**, Università Bocconi, Italy

Co-authors: Giacomo Zanella

Finding effective ways to exploit parallel computing to speed up MCMC convergence is an important problem in Bayesian computation and related disciplines. The zero-order (aka derivative-free) version of the problem is considered, where it is assumed that (a) the gradient of the target distribution is unavailable (either for theoretical, practical, or computational reasons) and (b) the (expensive) target distribution is evaluated in parallel at K different locations and these evaluations are used to speed up MCMC convergence. Two main contributions are made in this respect. First, it is shown that any method falling within a fairly general "multiple proposal framework" can only speed up convergence by $\log(K)$ factors in high dimensions. The fundamental limitation of such a framework, which includes multiple-try MCMC as well as many other previously proposed methods, is that it restricts possible moves to the support of the K evaluation points. Results are stated in terms of upper bounds on the spectral gap of the resulting scheme. Second, it is discussed how stochastic gradient estimators can be used to make better use of parallel computing and achieve polynomial speedups in K . Some of the methods have similarities, but also notable differences, with classical zero-order optimization methods.

C0758: **Sampling using time-changed Markov processes**

Presenter: **Giorgos Vasdekis**, Newcastle University, United Kingdom

Co-authors: Andrea Bertazzi

A framework of time-changed Markov processes is introduced to speed up the convergence of Markov chain Monte Carlo (MCMC) algorithms in the context of multimodal distributions and rare event simulation. The time-changed process is defined by adjusting the speed of time of a base process via a user-chosen, state-dependent function. This framework is applied to several Markov processes from the MCMC literature, such as Langevin diffusions and piecewise deterministic Markov processes, obtaining novel modifications of classical algorithms and also rediscovering known MCMC algorithms. Theoretical properties of the time-changed process are proven under suitable conditions on the base process, focusing on connecting the stationary distributions and qualitative convergence properties such as geometric and uniform ergodicity, as well as a functional central limit theorem. Time permitting, the approach will be compared with the framework of space transformations, clarifying the similarities between the approaches.

C0636: **Bayesian inference from time series of allele frequency data using exact simulation techniques**

Presenter: **Jaromir Sant**, Università di Torino, Italy

Co-authors: Paul Jenkins, Dario Spano, Jere Koskela

A central statistical problem in population genetics is to infer evolutionary and biological parameters such as the strength of natural selection and allele age from DNA samples extracted from a contemporary population. That all samples come only from the present-day has long been known to limit statistical inference; there is potentially more information available if one also has access to ancient DNA so that inference is based on a time-series of historical changes in allele frequencies. An MCMC method is introduced for Bayesian inference from allele frequency time-series data based on an underlying Wright-Fisher diffusion, through which one can infer the parameters of essentially any selection model, including those with frequency-dependent effects. The chief novelty is that the method is shown to be exact in the sense that it is possible to augment the state space with the unobserved diffusion trajectory, even though the transition function is intractable. Through careful design of a proposal distribution, we describe an efficient method in which updates to the trajectory and accept/reject decisions are calculated without error. The method is illustrated on data capturing changes in coat color over the past 20,000 years, and evidence is found to support previous findings that the mutant alleles ASIP and MC1R responsible for changes in coat color have experienced very strong, possibly overdominant, selection, and estimates are further provided for the ages of these genes.

C0632: **Learning with importance weighted variational inference**

Presenter: **Kamelia Daudel**, ESSEC Business School, France

Several popular variational bounds involving importance weighting ideas have been proposed to generalize and improve on the Evidence Lower Bound (ELBO) in the context of maximum likelihood optimization, such as the importance weighted auto-encoder (IWAE) and the variational Renyi (VR) bounds. The methodology to learn the parameters of interest with these bounds typically amounts to running stochastic gradient-based variational inference algorithms that incorporate the reparameterization trick. While the outcome of the resulting variational inference algorithms is expected to be tied to the choice of the variational bound, the precise effect of that choice remains poorly understood. The comparison of the ELBO, IWAE, and VR bounds methodologies is enabled by providing asymptotic analyses for key stochastic gradient estimators used in these methodologies. The analyses also reveal how these estimators compare to each other, and the theoretical findings are empirically illustrated.

C0359: **Extrapolation of tempered posteriors**

Presenter: **Mengxin Xi**, Kings College London, United Kingdom

Co-authors: Marina Riabiz, Chris Oates, Zheyang Shen, Nicolas Chopin

Accurately estimating quantities of interest from posterior distributions within a limited computational budget is essential across various fields. However, sampling from informative posterior distributions presents significant challenges. Tempering methods facilitate the construction of a path from the prior distribution to the complex posterior distribution. Instead of using samples from the posterior distribution to estimate posterior expectations, samples from intermediate tempered distributions and their corresponding tempered posterior expectations are used. The knowledge of the intermediate distributions enables posterior quantities of interest to be extrapolated. Specifically, weak sufficient conditions are established under which tempered expectations are not merely smooth as a function of t , but analytic, implying that knowledge of the tempered expectation in any open t interval fully determines the posterior expectation of interest.

C0044 Room BCB 213 ADVANCES IN QUANTITATIVE RISK MANAGEMENT AND COPULAS**Chair: Hideatsu Tsukahara****C0442: Modeling and forecasting the co-movement of international yield curve drivers****Presenter:** Maria Sprincenatu, Sofine, Meag Munich Ergo AssetManagement GmbH, Germany**Co-authors:** Stefan Mittnik

New data-driven state-space models are developed to forecast the co-movement of yield curve drivers of different world regions. The models are designed to preserve the dynamic properties of the yield curve drivers embodied in their underlying data generation processes. The models allow forecasting the co-movement of yield curves of different world regions by forecasting their drivers. Using actively traded government bond yields for the US and Germany, it is shown that the models outperform both the domestic and the global Diebold-Li models at long forecast horizons. It is also shown that the curvature does have predictive power for short forecast horizons.

C0524: Vine copula mixed models for network meta-analysis of multiple diagnostic tests**Presenter:** Aristidis Nikoloulopoulos, University of East Anglia, United Kingdom

As meta-analysis of multiple diagnostic tests impacts clinical decision making and patient health, there is growing interest in statistical models that synthesize evidence from studies comparing multiple diagnostic tests. To compare the accuracy of multiple diagnostic tests in a single study, three designs are commonly used: (i) the multiple test comparison design; (ii) the randomized design, and (iii) the non-comparative design. Generalized linear mixed models (GLMMs) are currently the recommended approach for jointly meta-analyzing data from all three designs, enabling simultaneous inference. In this context, vine copula mixed models are proposed as a flexible and powerful alternative. These models generalize the GLMM framework by allowing for arbitrary univariate distributions of the random effects and capturing tail dependencies and asymmetries. Findings indicate that vine copula mixed models can offer improvements over GLMMs, supporting their adoption for network meta-analysis of multiple diagnostic tests.

C0288: Model selection tests for truncated vine copulas under nested hypotheses**Presenter:** Ichiro Nishi, The Graduate University for Advanced Studies, Japan**Co-authors:** Yoshinori Kawasaki

Vine copulas provide a flexible framework for modeling multi-dimensional dependencies. However, this flexibility is accompanied by rapidly increasing complexity as dimensionality grows, necessitating appropriate truncation to manage this challenge. While the use of Vuong's model selection test has been proposed as a method to determine the optimal truncation level, its application to vine copulas has been heuristic, assuming only strictly non-nested hypotheses. This assumption conflicts with the inherent nesting within truncated vine copula structures. Vuong's model selection tests are systematically applied to distinguish competing models of truncated vine copulas under both nested and strictly non-nested hypotheses (Vuong-N and Vuong-SNN tests, respectively). Through extensive simulation studies, the performance of the Vuong-N and Vuong-SNN tests is evaluated using p-values, the number of rejections, and mean empirical KLIC. The results reveal that the relative performance of each test is sensitive to the strength of dependencies within the vine structure. In scenarios with weaker pairwise dependencies, the Vuong-N test produced lower p-values and higher rejection rates, along with improved mean empirical KLIC. Conversely, when the dependencies are stronger, the Vuong-SNN test yielded valid and often superior model distinctions, demonstrating that strictly non-nested testing, despite its heuristic status, remains an informative approach in such settings.

C0765: Forecasting stock returns using equi-correlation structures and component selection**Presenter:** Yoshinori Kawasaki, The Institute of Statistical Mathematics, Japan**Co-authors:** Takayuki Morimoto, Yohji Akama

A novel factor modeling approach is proposed for stock return prediction in the Japanese equity market by utilizing dynamic equi-correlation structures derived from daily industry returns. The dynamic equi-correlation model is implemented to estimate time-varying equi-correlation coefficients across TOPIX industry portfolios and construct an industry equi-correlation (IEC) index by removing medium-term trends. Applying principal component analysis to these IEC series, latent forecasting factors are extracted. To determine the number of components to retain, four statistically grounded selection rules are incorporated: Broken-stick rule, adjusted correlation thresholding, Guttman-Kaiser, and cumulative percentage of variance, based on recent theoretical contributions by prior studies. The predictive power of the resulting IEC-based components is evaluated using panel regressions, and model performance is assessed via out-of-sample R-squared, Sharpe ratio, and certainty equivalent return. Results show that the correlation-based factors significantly outperform traditional accounting-based factor models in terms of forecast accuracy and economic utility. This framework highlights the value of correlation dynamics and dimensionally efficient factor compression in building interpretable and robust asset pricing models tailored to the Japanese market.

C0829: Dynamic spatial panel data models with copulas and model validation**Presenter:** Hideatsu Tsukahara, Seijo University, Japan

The purpose is to consider a dynamic spatial panel data model with common factors as an extension of multifactor models in financial econometrics. It incorporates explanatory variables in the classical linear regression manner and utilizes a spatial weight matrix as a means to express network dependence. The underlying disturbance distribution is assumed to be Gaussian. By fitting a skew-t copula to the residuals, one can check the Gaussian assumption as well as whether the asymmetry in dependence structure and tail dependence have been captured by the variables already in the model. The aim is to propose estimation and testing methods for unknown parameters in the model, and investigate their properties. Some simulation results and empirical applications to real data are given.

C0293 Room BCB M201 ASYMPTOTICS, STATISTICS OF EXTREMES AND QUANTILE INFERENCE**Chair: Stefano Rizzelli****C0185: Asymptotically unbiased estimator of the extreme value index under random censoring****Presenter:** Armelle Guillou, Strasbourg University, France**Co-authors:** Martin Bladt, Yuri Goegebeur

The purpose is to consider bias-corrected estimation of the extreme value index of a Pareto-type distribution in the censoring framework. The initial estimator is based on a Kaplan-Meier integral from which the bias is removed under a second-order framework. This estimator depends on a suitable external estimation of second-order parameters, which is also discussed. The weak convergence of the bias-corrected estimator is established. It has the nice property of having the same asymptotic variance as the initial estimator. This nice feature is illustrated in a simulation study where the estimator is compared to alternatives already introduced in the literature. Finally, the methodology is applied to an insurance dataset.

C0253: A general theory for extremal regression in heavy-tailed models**Presenter:** Abdelaati Daouia, Toulouse School of Economics, France**Co-authors:** Yasser Abbas, Gilles Stupfler

Studying rare events at the heavy tails of conditional Pareto-type distributions, in the presence of high-dimensional covariates, is a burgeoning science with many applications in actuarial, financial, and environmental risk management. The most prominent risk measures to quantify these events utilize conditional quantiles, expected shortfall, and expectiles at extreme levels. The few attempts to tackle this extreme value problem involve location-scale regression models with heavy-tailed noise. A more flexible and complex model is employed that better balances model generality with estimation efficiency. A general theory is developed that relies on residual-based estimators of the three regression risk measures at

both intermediate and extreme levels, and their asymptotic behavior is fully explored in generic settings. Simple sufficient criteria are also provided for verifying the main high-level assumption, which facilitates the construction of weighted Gaussian approximations for the tail quantile residual process, ultimately ensuring the asymptotic normality of all produced extreme value estimators. This generic extremal regression framework is then applied to linear, nonlinear, and nonparametric estimation scenarios. Simulations show the undeniable potential of the methodology for various distribution types, outperforming the best available competing estimation approaches. An application to real financial data further solidifies their dominance.

C0638: **Asymptotics of tail pairwise dependence matrices**

Presenter: **Stephane Lhaut**, ENSAE Paris, France

Co-authors: Johan Segers, Anna Kiriliouk

In recent literature on multivariate extremes, various definitions of tail pairwise dependence matrices (TPDMs) have been proposed to summarize the tail dependence structure of a random vector X . Most of these approaches rely on the assumption that X is multivariate regularly varying, which implies that all marginal tails are equivalent. Several popular definitions of TPDMs are unified, and it is shown how nonparametric inference can be carried out in each case without requiring multivariate regular variation of X . Instead, regular variation of a standardized version V is assumed, where all margins have been transformed to a common scale, a more realistic assumption in many applications. The joint asymptotic normality of TPDM entries is established, based either on the empirical stable tail dependence function or on the empirical angular measure, depending on the setting. Applications include inference for parametric models and dimension reduction via multidimensional scaling.

C0679: **Extremes extrapolation in time series: Accurate Bayesian inference based on the peaks over threshold method**

Presenter: **David Carl**, Bocconi University, Italy

Co-authors: Simone Padoan, Stefano Rizzelli

A strictly stationary time series is considered with marginal distribution in the domain of attraction of the generalized extreme value (GEV) distribution. Under mild conditions concerning the serial dependence structure, the largest observations after a linear transformation, i.e., the normalized peaks over a threshold, converge in distribution to a generalized Pareto (GP) distribution. This motivates likelihood-based inference using the GP distribution. It is shown that the resulting naive Bayesian approach, treating the observations as independent, will lead to credible intervals that fail to achieve asymptotic correct coverage. An adjustment is proposed to the likelihood to remedy this issue and show that the resulting posterior distributions for the GP parameters and extreme quantiles retain the same contraction rates as under independence while simultaneously allowing for accurate uncertainty estimation. If there is time left, it will be explained how to extend these results in order to achieve dynamic extrapolation of future extremes.

C1015: **Fast and efficient inference for flexible spatial extremes models**

Presenter: **Boris Beranger**, University of New South Wales, Australia

Co-authors: Scott Sisson, Peng Zhong

Statistical modelling of spatial extreme events has gained increasing attention over the last few decades, with max-stable processes and, more recently, r-Pareto processes becoming the reference tools for the statistical analysis of asymptotically dependent data. Although inference for r-Pareto processes is easier than for max-stable processes, there remain major hurdles for their application to very high-dimensional datasets within a reasonable timeframe. In addition, both approaches have almost exclusively considered the Brown-Resnick model for its Gaussian-based foundations and the continuity of its exponent measure. A class of models is derived, for which this continuity property holds, and the skewed Brown-Resnick model is presented, an extension of the Brown-Resnick that allows for non-stationarity in the dependence structure, and the truncated extremal-t, a refinement of the well-known extremal-t model. An inference methodology is used based on the intensity function of the process, which is derived from the exponent measure, and the statistical and computational efficiency of this approach is demonstrated. Applications to two real-world problems illustrate valuable gains in flexibility from the proposed models as well as appealing computational gains over reference methodologies.

CO415 Room BCB M202 CLIMATE FINANCE

Chair: Julien Chevallier

C0380: **Green commodity futures investing with style**

Presenter: **Joelle Miffre**, Audencia, France

Co-authors: Mascia Bedendo, Adrian Fernandez-Perez, Ana-Maria Fuertes

This paper introduces a sustainable, style-integrated commodity portfolio designed to balance strong financial returns, leveraging cross-sectional return predictors with improved environmental outcomes by minimizing exposure to high-emission (brown) commodity futures. While the sustainable style-integrated portfolio underperforms a purely style-based portfolio in financial terms, it offers greater decarbonization benefits and enhanced resilience to climate-related risk shocks. The robustness of the framework is demonstrated across various emissions metrics, integration strategies, and tilting functions. Constrained capital flows into brown commodity futures post the 2015 Paris Agreement and hedging pressure are identified as the economic channels through which portfolio decarbonization translates into climate change mitigation.

C0708: **Biodiversity risk in commodity futures**

Presenter: **Zheng Zhang**, City University London, United Kingdom

Co-authors: Ana-Maria Fuertes, Kate Phylaktis

The purpose is to investigate whether biodiversity risk is priced in commodity futures markets and whether such risk contributes to downside tail exposure. While prior studies have focused on climate risk, we examine biodiversity as a distinct environmental risk. The analysis uses a news-based biodiversity risk index from The New York Times constructed at the daily frequency by a prior study, and evaluates 25 commodity futures across agriculture, energy, and metals. A multi-method approach is adopted, combining Fama-MacBeth regressions, lag-augmented local projections, and event studies. Downside risk is also estimated using value-at-risk and second-order lower partial moments. Findings suggest that biodiversity risk is not persistently priced in the cross-section of commodity returns. However, significant but short-lived abnormal returns are observed around biodiversity-related policy shocks, particularly under heightened macro-policy uncertainty. Furthermore, unexpected biodiversity shocks contribute meaningfully to downside risk, especially in agriculture and energy markets. These results suggest that biodiversity risk is not yet a systemic factor in commodity pricing but is becoming increasingly relevant in event-driven, state-contingent contexts. The findings highlight the need for improved biodiversity risk disclosure and underscore the importance of incorporating biodiversity scenarios into risk management and stress testing frameworks.

C0370: **Investor sentiment and green shipping: Combining proxy variables and a large language model**

Presenter: **Ioannis Moutzouris**, City, University of London, United Kingdom

Co-authors: Yao Shi, Nikos Papapostolou, Panos Poulisias

Existing research suggests that investor sentiment affects asset prices, returns, and future investment. Nonetheless, little is known about its influence on investment in green assets. The example of a real-asset, capital-intensive industry (shipping) is taken to study the effect of sentiment on green investment. The sentiment indices are built with large language models, bidirectional encoder representations from transformers (BERT) and financial BERT (FinBERT), as well as key shipping finance variables. The focus is on the period 2020-2025, where the regulations for the reduction of greenhouse gas emissions for shipping become increasingly strict. The results indicate that positive news on the industry's net-zero transition

significantly negatively influences future investment in greener vessels. This contrasts with the documented evidence that positive sentiment drives investment in conventional (i.e., non-green) assets. Findings suggest that an overly optimistic outlook on the net-zero progress may hinder green initiatives, likely because investors perceive less urgency to act. Findings yield important implications for the industry and policymakers alike.

C0346: Contagion vs competition: Evidence from crypto exchange flows

Presenter: Kirill Shakhnov, University Of Surrey, United Kingdom

Co-authors: Daniele Bianchi

The purpose is to investigate how shocks affecting cryptocurrency exchanges propagate through the network of Bitcoin flows. It distinguishes between two channels of propagation: Contagion and competition. Using a novel dataset of on-chain Bitcoin flows between exchanges, it is analyzed how exchange-specific shocks affect network flows, liquidity, and price discovery. It is found that negative shocks to exchanges lead to significant changes in inter-exchange flows and affect market quality measures, including price differences and bid-ask spreads. Evidence suggests that the cryptocurrency market structure exhibits both competitive dynamics and contagion effects, with the relative importance depending on the nature of the shock and the specific exchange affected.

C0688: Uncovering climate transition risk in FX markets through equity risk premia

Presenter: Ana-Maria Fuertes, City University London, United Kingdom

Co-authors: Kate Phylaktis, Zheng Zhang

The aim is to investigate whether climate transition risk is priced in foreign exchange (FX) markets and through which channels it transmits. Using monthly data for 35 countries (15 developed, 20 emerging) from 2000 to 2024, it is examined how climate transition risk affects FX returns via equity market risk premia. A dual empirical approach is adopted, combining portfolio-sorting and asset pricing models. Climate transition risk is proxied by a country-level index based on greenhouse gas efficiency and emission pledges. It is found that climate transition risk is significantly priced in emerging markets, especially after the 2015 Paris Agreement. High-risk emerging economies exhibit elevated equity risk premia and negative FX returns, reflecting systematic deviations from uncovered equity parity (UEP). These results are robust across methods, investment horizons, and alternative risk proxies. In contrast, developed markets show no consistent pricing of transition risk. Findings suggest that international investors rebalance portfolios in response to climate policy uncertainty, particularly in less-regulated, more vulnerable markets. This has implications for global capital flows, FX exposure management, and sovereign financing in the low-carbon transition. The results underscore the growing importance of integrating climate risk into currency risk management, especially for policymakers and investors active in emerging markets.

C1367: Spectral climate risk

Presenter: Andrea Cipollini, University of Palermo, Italy

Co-authors: Iolanda Lo Cascio, Fabio Parla, Fabio Parla

The purpose is to examine the return performance of a green-minus-brown (GMB) portfolio aiming to hedge climate risk for US. While existing studies focus on the empirical analysis within the time domain, the aim is to analyze the contribution of climate risk to return performance that varies across frequency bands. For this purpose, the extended Wold decomposition of risk-adjusted return and climate concern shock is used. Then, in a second stage, regression analysis shows that green stocks outperform brown stocks over short-medium term horizon.

CO153 Room BCB 308 STATISTICAL MODELING OF TEXT DATA

Chair: Bettina Gruen

C0268: Nested Dirichlet models for unsupervised attack pattern detection in honeypot data

Presenter: Francesco Sanna Passino, Imperial College London, United Kingdom

Co-authors: Anastasia Mantziou, Daniyar Ghani, Philip Thiede, Ross Bevington, Nick Heard

Cyber-systems are under near-constant threat from intrusion attempts. Attack types vary, but each attempt typically has a specific underlying intent, and the perpetrators are typically groups of individuals with similar objectives. Clustering attacks that appear to share a common intent is very valuable to threat-hunting experts. Dirichlet distribution topic models are explored for clustering terminal session commands collected from honeypots, which are special network hosts designed to entice malicious attackers. The main practical implications of clustering the sessions are two-fold: Finding similar groups of attacks and identifying outliers. A range of statistical models is considered, adapted to the structures of command-line syntax. In particular, concepts of primary and secondary topics, and then session-level and command-level topics, are introduced into the models to improve interpretability. The proposed methods are further extended in a Bayesian nonparametric fashion to allow unboundedness in the vocabulary size and the number of latent intents. The methods are shown to discover an unusual MIRAI variant that attempts to take over existing cryptocurrency coin-mining infrastructure, not detected by traditional topic-modeling approaches.

C0296: Narrative shift detection: A hybrid approach of dynamic topic models and large language models

Presenter: Kai-Robin Lange, TU Dortmund University, Germany

As narratives in media evolve rapidly, understanding and investigating how narratives develop over time has become increasingly important. While large language models (LLMs) are effective at capturing narrative elements, small research groups might struggle to apply them across entire corpora due to high computational and financial costs. To address this issue, a method is introduced that combines the language understanding capabilities of LLMs with the scalability of dynamic topic models to analyze narrative shifts over time, utilizing the narrative policy framework. A dynamic topic model, along with a change point detection algorithm, is used to identify topical changes. Documents representative of these changes are then selected and analyzed using an LLM to automatically interpret the nature of the change and distinguish between narrative shifts and mere content shifts. This approach is applied to a corpus of "The Wall Street Journal" articles spanning 2009 to 2023. Results suggest that LLMs are effective at extracting narrative shifts when such shifts are present, but are limited when distinguishing between content changes and genuine narrative transitions.

C0337: Sampling uncertainty of research topics

Presenter: Anna Staszewska-Bystrova, University of Lodz, Poland

Co-authors: Victor Bystrov, Viktoriia Naboka-Krell, Peter Winker

In latent topic models, estimated topic-word and document-topic probabilities are typically reported with no indication of sampling uncertainty. The lack of additional information on sampling uncertainty might result in misleading conclusions regarding topic structure and prevalence. The proposal is to measure sampling uncertainty using a bootstrap method and describe how uncertainty can be captured by novel types of word clouds reporting topic-word probability estimates and by confidence bands designed for reporting time series estimates of topic weights. The application of the new measures and methods is illustrated with an empirical example involving conference abstracts. The results indicate varying robustness of estimated research topics with respect to resampling of documents from the same text collection. In particular, some estimated topics may not persist across resampled corpora, and the estimation precision of topic-word probabilities within the same topic can exhibit significant variation. Similar uncertainty is associated with topic prevalence over time. The proposed confidence bands for dynamic topic weights can be used to make inferences about structural changes in research topic trends.

C0336: Seeded Poisson factorization topic models with covariates

Presenter: Bettina Gruen, WU Vienna University of Economics and Business, Austria

Co-authors: Bernd Probstmaier, Paul Hofmarcher

Topic models infer latent structures in text corpora to guide data-driven detection of themes and to cluster or group documents. The basic topic model only requires the transformation of the documents in the corpus into a document-term matrix to perform inference based on either the latent Dirichlet allocation or the Poisson factorization models. Many extensions have been proposed and considered to improve the insights gained in applications and allow for the inclusion of additional information, such as the inclusion of seed words to guide topic discovery or covariates to infer, for example, associations between document characteristics and topic distributions. Many of these extensions build on the latent Dirichlet allocation model, with, for example, keyATM including seed words and the structural topic model allowing for covariates to be taken into account. The focus is on the Poisson factorization model, and seeded Poisson factorization is extended to include covariates that drive topic distributions of documents. The estimation is investigated using variational inference to allow for large-scale performance and empirically assess the performance of applying the model, including tools for suitable post-processing and model inspection.

C0354: Framing the evidence: A text mining approach to regulatory persuasion in pharmaceutical pricing

Presenter: **Paul Hofmarcher**, University Salzburg, Austria

Pharmaceutical companies engage in high-stakes persuasion when submitting benefit assessment dossiers to health authorities. In Germany, this process requires firms to demonstrate the added value of new drugs to negotiate reimbursement prices. While the scientific evidence provided in these submissions is standardized, it is hypothesized that the "framing of information" may influence the reimbursement prices. A novel statistical framework, the structural text-based scaling (STBS) model, is proposed, which combines Poisson factorization topic modeling with author-level framing effects to detect variation in textual emphasis across those dossiers. The Bayesian hierarchical model allows topic-specific deviations in language use to be regressed on covariates such as company identity, orphan drug status, and therapeutic class. Estimation is performed via variational inference. Using a corpus of oncological drug dossiers submitted to regulatory authorities, preliminary results suggest modest but systematic differences in linguistic framing across firms and drug categories. This might have potential implications for regulatory pricing decisions.

CO185 Room BCB 309 HIGH-DIMENSIONAL DATA ANALYSIS

Chair: Marcell Tamas Kurbucz

C0321: Bias-variance trade-off in feature selection for generalized additive models under concurvity

Presenter: **Laszlo Kovacs**, Corvinus University of Budapest, Hungary

Co-authors: Tibor Keresztely, Zoltan Madari

The bias-variance trade-off in statistical learning is about finding a balance between two types of errors: Bias, which is the difference between the expected value of an estimate and the population parameter, and variance, which is the variability of estimates around their expected value. In feature selection, this trade-off plays a critical role in model performance. Selecting too few features may lead to high bias, as the model becomes overly simplistic and may fail to capture important patterns. Conversely, selecting too many features can result in high variance, as the model may overfit to noise in the training set. In generalized additive models (GAMs), concurvity - a non-linear extension of multicollinearity - causes further inflation in the variance of estimators. Some feature selection algorithms attempt to address concurvity, but only for the pairwise cases, so they do not consider when a feature is a multivariate function of several other variables. GAM feature selection algorithms are compared in Monte Carlo simulations under different high-dimensional concurvity scenario groups: No concurvity, pairwise concurvity and multivariate concurvity. The object of the comparisons is the bias-variance trade-off in the spline estimates of GAMs proposed by the examined feature selection algorithms.

C0328: SplitWise regression: Stepwise modeling with adaptive dummy encoding

Presenter: **Marcell Tamas Kurbucz**, University College London, United Kingdom

Co-authors: Nikolaos Tzivanakis, Nilufer Sari Aslam, Adam Sykulski

Capturing nonlinear effects while preserving interpretability remains a key challenge in regression modeling. SplitWise, a novel extension of stepwise regression, introduces adaptive transformations of numeric predictors into binary threshold-based features using shallow decision trees. These transformations are only applied when they improve model fit based on Akaike or Bayesian information criteria (AIC/BIC). This approach enhances the flexibility of linear models without sacrificing their transparency. Implemented as an R package, SplitWise is validated on synthetic and real-world datasets. Compared to traditional stepwise and penalized regression, it consistently produces more parsimonious, interpretable, and generalizable models.

C0967: Supervised learning of binary features by decision trees for weightless neural networks

Presenter: **Douglas Cardoso**, University of Coimbra, Portugal

An original method based on decision trees for the supervised learning of binary features is discussed. As a case study, it was used to improve the performance of weightless neural networks (WNNs), whose most usual mathematical modeling presumes binary inputs, unlike other well-known neural networks theorized to directly operate on real-valued data. A systematic methodology is introduced for the extraction of binary features through decision tree learning, optimizing for both classification accuracy and feature sparsity, leveraging the power of decision trees to efficiently identify discriminative binary features from input data. Experimental evaluation on various classification datasets from the OpenML CC18 benchmark suite confirms the effectiveness of the approach in generating compact yet discriminative representations, which provided a statistically significant accuracy gain compared to the arguably most used binarization method for WiSARD, a reference model of WNN.

C1313: Causality with or without the Reichenbach principle: The degrees of freedom method

Presenter: **Andras Telcs**, HUN-REN Wigner RCP, Hungary

The purpose is to introduce a novel methodology for uncovering causal interactions between deterministic or stochastic dynamic systems, with a particular focus on distinguishing true causal relationships from misleading correlations arising due to hidden common causes. The proposed approach is applicable to both discrete and continuous systems in time and space. Central to the method is random forest clustering of coupled 0-1 Markov chains, through which causal relationships between systems are classified directly from observed data. The effectiveness of the technique is demonstrated on both simulated and real-world time series. Notably, the method reveals a new result: the number of sunspots influences global Earth temperature with zero lag on a monthly scale. To the best of knowledge, prior studies have reported minimal lags of two to three years, with some suggesting even longer intervals.

C1262: Representation of reality

Presenter: **Antal Jakovac**, Wigner Research Centre for Physics, Hungary

Present-day AI applications are often based on stateless models, in which the representation of the actual reality is much simpler than that of general knowledge. A typical example is large language models (LLMs), where the actual reality is encoded in the conversation history, as opposed to the hundreds of billions of parameters used to represent general knowledge. Retrieval-augmented generation (RAG) systems are being developed to overcome this limitation, but considerable effort is still needed to achieve satisfactory solutions. The aim is to overview the theoretical foundations of representing reality and propose a data storage logic that reflects these ideas.

CO025 Room BCB 310 STATISTICAL METHODS AND THEIR APPLICATIONS

Chair: Lorenzo Mercuri

C0322: On some mobility measures for Markov processes with application to a wind-power production

Presenter: **Guglielmo Damico**, University G. d'Annunzio of Chieti-Pescara, Italy

The rate of occurrence of failures (ROCOF) is a standard metric for evaluating system performance over time, but it falls short in capturing a

system's instantaneous behavior. Three new complementary indicators are introduced: Rate of occurrence of repairs (ROCOR), rate of inoccurrence (ROI), and total mobility rate (TMR), to provide a richer, time-dependent view of system reliability, especially in Markov systems. ROCOR quantifies the system's immediate tendency to transition from failure to working states. ROI measures the likelihood of remaining within the current subset of states (working or failed). TMR integrates ROCOF, ROCOR, and ROI to reflect the system's overall dynamism. Explicit formulas are derived for these metrics in the context of Markov models. Their practical value is demonstrated through a real-world application to wind farm management, where these indicators reveal nuanced operational differences between sites that share similar long-term wind profiles.

C0364: Is Spearman's a robust measure of correlation for financial time series? Some evidences in the insurance sector

Presenter: **Roberto Baviera**, Politecnico di Milano, Italy

Spearman's rho is a widely used non-parametric measure of statistical dependence between two time series. Since it is rank-based, it is considered robust, as it is less sensitive to extreme values and therefore relatively resistant to outliers. It is proven that when applied to financial time series with a significant number of zeros, Spearman's rho can produce widely varying correlation estimates. In such cases, financially equivalent time series - those that should theoretically exhibit similar dependence - can record highly different Spearman's correlation. An experimental analysis is performed using datasets in the insurance sector.

C0562: A new model for the perceived time to transition to a low carbon economy

Presenter: **Edit Rroji**, Università degli studi di Milano-Bicocca, Italy

Co-authors: Lorenzo Mercuri, Andrea Perchiazzo, Ilaria Stefani

In the context of the transition to a low or zero carbon economy, the difference in greenium between pairs of twin bonds with different maturities is expected to disappear or, at least, to reduce in both level and volatility. Consequently, a model is needed that imposes a terminal condition on the dynamics of the process representing the difference in nodes within the greenium term structure. An important feature of this difference, observed in empirical data, is its mean-reverting behavior. This characteristic motivates the introduction of ad-hoc models that consider the possibility of a transition occurring at a specific time. Two models are discussed: The first is an extension of the classical Vasicek model, where the volatility term remains constant until a future time instant, after which it decreases linearly. This model is integrated into a regime-switching framework, where the perceived deadline for transitioning to a low or zero-carbon economy defines the regime. Both models are calibrated using market data extracted from twin German government bonds.

C0239: When Tukey meets Chauvenet: A new boxplot criterion for outlier detection

Presenter: **Hongmei Lin**, Shanghai University of International Business and Economics, China

The box-and-whisker plot is one of the most popular graphical methods in descriptive statistics. On the other hand, however, Tukey's boxplot is free of sample size, yielding the so-called "one-size-fits-all" fences for outlier detection. Although improvements on the sample size adjusted boxplots do exist in the literature, most of them are either not easy to implement or lack justification. As another common rule for outlier detection, Chauvenet's criterion uses the sample mean and standard deviation to perform the test, but it is often sensitive to the included outliers and hence is not robust. By combining Tukey's boxplot and Chauvenet's criterion, a new boxplot, namely the Chauvenet-type boxplot, is introduced with the fence coefficient determined by an exact control of the outside rate per observation. The new outlier criterion not only maintains the simplicity of the boxplot from a practical perspective, but also serves as a robust Chauvenet's criterion. Simulation study and a real data analysis on the civil service pay adjustment in Hong Kong demonstrate that the Chauvenet-type boxplot performs extremely well regardless of the sample size, and can therefore be highly recommended for practical use to replace both Tukey's boxplot and Chauvenet's criterion.

C0294: Examining directional association between depression and anxiety

Presenter: **Soumik Purkayastha**, University of Pittsburgh, United States

Utilizing a novel entropy loss (EL) metric, a causal discovery method is proposed to understand directional effects in the causal relationship between depression and anxiety among medical interns. This method advances existing methods of bivariate causal discovery with theoretical guarantees of causal effect identifiability and statistical inference, and enjoys good computational performance. Using data from the Intern Health Study (n=6,858), the proposed method reveals with high statistical confidence that depression scores (PHQ-9) consistently predispose anxiety scores (GAD-7) across four longitudinal visits, controlling for demographic confounders. This finding provides crucial insights into the directional effect, useful for mental health intervention strategies for medical interns. Simulation studies demonstrate that EL achieves nearly superior accuracy compared to existing approaches across various conditions with reduced computation time. The EL framework's ability to handle discrete clinical scores while adjusting for confounders makes it particularly valuable for psychiatric epidemiology and broader applications in causal discovery with discrete data.

CO141 Room BCB 311 COMPLEX DATA UNCOVERED: METHODOLOGY AND PRACTICAL INSIGHTS

Chair: Elena Ballante

C0666: A novel cluster-weighted multilevel model for two-levels clustering

Presenter: **Chiara Masci**, Università degli Studi di Milano, Italy

Co-authors: Andrea Cappozzo

A novel two-level cluster-weighted model designed for hierarchical data is introduced. This approach bridges cluster-weighted modeling with multilevel modeling using discrete random effects, resulting in a powerful framework for two-level clustering. The model identifies latent subpopulations of observations - referred to as profiles - that exhibit heterogeneous characteristics. Within each profile, multilevel models are estimated to capture profile-specific regressions and random effects. At the higher level of the hierarchy, random effects are modeled using a semi-parametric generalized linear mixed model (SPGLMM). The discrete nature of these random effects facilitates the identification of clusters at the top level of the data structure. To estimate the model parameters, an expectation-maximization algorithm is developed, tailored to the proposed framework. The method is validated through a simulation study and an empirical case study. It is particularly well-suited to applications involving nested data structures, such as students within schools, citizens within states, or patients within hospitals. To demonstrate its practical utility, the model is applied to the OECD-PISA dataset to analyze European educational systems. Specifically, the probability of a student being a low performer in mathematics is examined, identifying distinct student profiles and grouping countries according to their impact on student outcomes.

C1012: Clustering longitudinal data in clinical studies: A practical comparative analysis of statistical methods

Presenter: **Paola Rancoita**, Vita-Salute San Raffaele University, Italy

Co-authors: Nicolo Pecorelli, Chiara Brombin

Frequently, in clinical studies, patient status is monitored over time to decide patient management. These longitudinal data are challenging to analyze due to their intrinsic nature, irregular timing of measurements, missing values, and eventual nonlinearities. Moreover, often multivariate modeling (which jointly evaluates multiple outcomes) is needed to capture the multidimensional nature of the phenomenon. In this context, subject-specific heterogeneity frequently emerges, thus requiring the usage of robust statistical methods to identify latent patient subgroups, possibly related to their characteristics. This is particularly relevant for clinicians, as identifying clusters can allow the definition of personalized care strategies. From a statistical perspective, several clustering approaches for longitudinal data have been proposed, including algorithm-based (e.g., hierarchical clustering) and model-based approaches (e.g., finite mixture models, latent class mixed models). These methods differ in the number of cluster selection criteria, ability to include covariates or nested random effects, and to handle timing irregularities or complex trajectories. The flexibility

and performance of these clustering strategies in uncovering latent patient longitudinal profiles will be compared. Their advantages and limitations are shown by applying them to real longitudinal data assessing quality of life and functional capacity in patients after pancreatic resection.

C1039: Explainable ensemble clustering through mutual information, with applications on high dimensional data

Presenter: **Federico Maria Quetti**, University of Pavia, Italy

Co-authors: Elena Ballante, Paolo Giudici, Silvia Figini

Unsupervised learning techniques aim to uncover the intrinsic structure of data, with clustering being the process of grouping similar points together. A common limitation of many machine learning tasks is the lack of explainability of the process, which often operates as a black box. In clustering settings, a major challenge for most methods is the limited interpretability, as little insight is provided into which features drive the grouping, especially in high-dimensional settings. To address this limitation, a bagging-based clustering approach incorporating feature dropout is proposed, analogous to the supervised random forest methodology, aimed at decorrelating features in the partitioning steps. The involvement in clustering of each feature is ranked using an index based on information theory. In each step, the mutual information $I(X;Y) = H(X) - H(X|Y)$ ($H(X)$ being the Shannon entropy, $H(X|Y)$ the conditional entropy) between each feature involved and the estimated label obtained by the partitioning algorithm is evaluated. Then, an aggregated estimate is produced, weighing each step's contribution by an index of validity of the clustering (e.g., Dunn, Silhouette) to emphasize well-formed partitions. Results are presented on simulated and real datasets, with applications in medicine.

C1041: High-dimensional MANOVA test for semicontinuous biomedical data: Methodology and applications

Presenter: **Elena Sabbioni**, University of Oxford, United Kingdom

Co-authors: Claudio Agostinelli, Alessio Farcomeni, Elena Sabbioni

Modern biological applications increasingly involve complex data structures, driven by advances in technology that allow for the collection of a growing number of high-resolution features. In this context, the focus is on semicontinuous high-dimensional data, a common data type found in various fields, including genetics and medical research. These data combine a continuous part with positive observations and a part with exactly zero components, and often feature more variables than observations. This setting poses challenges for standard statistical methods, which typically address either semicontinuous or high-dimensional data, but not both simultaneously. To address this methodological gap, a novel MANOVA testing procedure is proposed, specifically designed to handle both of these features simultaneously. The method is based on a regularized likelihood ratio test, where the form of the penalized likelihood enables closed-form estimators, ensuring both computational tractability and scalability. Since the null distribution of the test statistic is unknown in this framework, we employ a permutation scheme. The efficiency of the method, both in terms of level and power of the test, is achieved in a simulation setting. Finally, its practical utility is illustrated through the analysis of a real microRNA expression dataset, showcasing its relevance for complex biomedical data.

C1076: Uncovering latent molecular patterns in mass spectrometry imaging via spatially-constrained graphical mixture models

Presenter: **Giulia Capitoli**, University of Milano-Bicocca, Italy

Mass spectrometry is a class of imaging techniques that measure molecular abundance in tissue samples at cellular resolution, while preserving the spatial structure of the tissue. In particular, mass spectrometry imaging has the capability to differentiate regions that are indistinguishable to pathologists at the microscopic level. A central goal in mass spectrometry data analysis is to identify molecules with similar functions within the analyzed biological system, enabling a better understanding of abnormal molecular mechanisms. The aim is to identify relevant biomolecules associated with cancer cells and the tumor microenvironment, thereby expanding biological knowledge. A Gaussian graphical mixture model is proposed to address unobserved heterogeneity and segment tissue sections into regions based on distinct molecular profiles. Specifically, the aim is to identify groups of molecules with similar activation patterns, investigate their spatial mapping within cancer tissue samples (e.g., renal and/or thyroid neoplasm), and discover clusters of molecules whose activation is linked to specific biological mechanisms. To model this heterogeneity, underlying molecular graphs are reconstructed from the data using sparsity constraints and spatial dependencies between neighboring pixels are incorporated. To account for the spatial nature of the dataset, hidden Markov random fields are utilized, ensuring that the spatial structure is effectively captured.

CO053 Room BCB 403 LATENT STRUCTURE IN COMPLEX DATA: BAYESIAN AND FREQUENTIST VIEWS

Chair: Silvia Montagna

C0669: Long memory network time series

Presenter: **Chiara Boetti**, University of Bath, United Kingdom

Co-authors: Matthew Nunes, Marina Knight

Many scientific areas, from computer science to the environmental sciences and finance, give rise to multivariate time series that exhibit long memory. Efficient modeling and estimation in such settings is key for a number of analysis tasks, such as accurate prediction. However, traditional approaches for modeling such data, for example, long memory vector autoregressive processes, are challenging even in modest dimensions, as the number of parameters grows quadratically with the number of modeled variables. In many data settings, the observed series is accompanied by a (possibly inferred) network that provides information about the presence or absence of between-component associations. Two new models are proposed for capturing the dynamics of long memory time series, where a network is taken into consideration. The approach not only facilitates the analysis of graph-structured long memory time series, but also improves computational efficiency over traditional multivariate long memory models by leveraging the inherent low-dimensional parameter space. Likelihood-based estimation algorithms are adapted to the network setting. Simulation studies show that the parameter estimation is more stable than traditional models and is able to tackle data scenarios where these models fail due to computational challenges. The efficacy of the proposed models is demonstrated on datasets arising in environmental science and finance applications.

C0635: Comparing flexible modelling approaches: The varying-thresholds model versus quantile regression

Presenter: **Marta Pittavino**, Ca Foscari University of Venice, Italy

Co-authors: Niccolo Ducci, Leonardo Grilli

The varying-thresholds model (VTM) is a novel methodology proposed by a prior study, capable of estimating the whole conditional distribution of a response variable in a regression setting. It can be used for continuous, ordinal, and count responses. The focus is on conditional quantiles and prediction intervals estimated through VTM, which are compared with those of quantile regression. The comparison is based on a set of data-generating models to assess the performance of the two methodologies regarding the coverage and width of prediction intervals. The simulation study encompasses settings with several functional forms and types of errors. In addition, a discrete version of the continuous ranked probability score is proposed as a tool to choose the best link function for the binary models used in the fitting of VTM. In summary, the varying-thresholds model is a flexible methodology that can be broadly applied with light assumptions; it is advantageous over quantile regression when the conditional quantile function is misspecified.

C1114: A Bayesian evidence synthesis model for estimating Hepatitis C prevalence among people who inject drugs in England

Presenter: **Pantelis Samartsidis**, University of Cambridge, United Kingdom

Co-authors: Daniela De Angelis, Matthew Hickman

Hepatitis C (HCV) is a blood-borne virus affecting the liver and can lead to acute liver damage, cirrhosis, and death. The introduction, since 2015, of effective drugs has led many countries to commit to elimination. Estimating prevalence among people who inject drugs, the most affected by HCV, is crucial to monitor progress towards elimination, but poses challenges. First, information on prevalence exists in many data sources, including surveillance and bio-behavioral surveys. The datasets are heterogeneous in terms of the populations that they target and the sample

size. Second, prevalence trends post 2015 differ by geographical region due to varying levels of treatment and by risk group due to differences in healthcare engagement. Third, in many data sources, behavioral characteristics are not recorded, hindering risk group classification. Finally, there is little information about the relative proportions of risk groups in the population, which is required to evaluate overall prevalence. To address these issues, a Bayesian evidence synthesis model is proposed that estimates prevalence by year, region, and risk group. A full characterization of uncertainty is obtained by accounting for potential misclassification of individuals in risk groups, and by modelling the relative proportions of risk groups using survey data. To tackle computational concerns, an integrated nested Laplace approximation is developed within Gibbs algorithm. The method is applied to data from England.

C0341: **Modeling count compositions through a structured mixture of Dirichlet-multinomial components**

Presenter: **Roberto Ascari**, University of Milano-Bicocca, Italy

Count compositions, or vectors of non-negative integers summing to a fixed total, are typically modeled using the multinomial distribution. While widely used, the multinomial has some limitations, particularly in its ability to capture positive covariance structures. One strategy to address this involves compounding the multinomial with a distribution defined on the simplex. A well-known example is the Dirichlet-multinomial (DM), which adds a parameter and improves the fit to real-world data, though it still imposes quite strong constraints on the covariance matrix. A novel distribution is introduced for count compositions, derived by compounding the multinomial with the extended flexible Dirichlet distribution. The resulting model can be seen as a structured finite mixture of specific DM components, allowing for greater flexibility and interpretability through latent group structures. Interestingly, it also allows for positive covariances. A regression model is developed based on this distribution, and its performance is evaluated through simulations and a real data application. Inference is conducted using a Bayesian framework, implemented via the Hamiltonian Monte Carlo algorithm.

C0593: **Modeling Eurovision 2025 songs with a flexible topic modeling approach**

Presenter: **Alice Giampino**, University of Milano-Bicocca, Italy

Co-authors: Roberto Ascari, Sonia Migliorati

Understanding the main topics of song lyrics in large-scale music competitions, such as the Eurovision Song Contest, is key to analyzing cultural trends, audience resonance, and artistic strategies. In 2025, with an increasing number of participating entries written in English, automatic methods for identifying underlying themes are especially relevant. However, lyrical data present significant challenges due to brevity, stylistic variability, and lexical diversity. Latent Dirichlet Allocation (LDA) has been widely used for uncovering latent topics by modeling word co-occurrence patterns. LDA's standard Dirichlet prior imposes only negative dependence among topics, limiting its ability to reflect nuanced thematic relationships or overlapping emotional tones in song lyrics. To overcome these constraints, a generalization of the LDA model is proposed through an extended flexible Dirichlet mixture prior tailored for topic distribution in musical texts. This enriched prior allows for positive correlations among topics, capturing common thematic clusters such as empowerment or political commentary that frequently co-occur across entries. The model maintains computational efficiency via conjugacy with the multinomial likelihood, enabling scalable inference through a collapsed Gibbs sampler. This approach yields a more expressive representation of lyrical themes, offering deeper insight into the collective voice of Eurovision 2025's songs and the cultural narratives they convey.

CO239 Room BCB 405 MODEL-BASED CLUSTERING: THEORY AND APPLICATIONS

Chair: Gertraud Malsiner-Walli

C0340: **Bayesian cluster weighted Gaussian models**

Presenter: **Panagiotis Papastamoulis**, Athens University of Economics and Business, Greece

Co-authors: Konstantinos Perrakis

A novel class of Bayesian mixtures is introduced for normal linear regression models, which incorporates a further Gaussian random component for the distribution of the predictor variables. The proposed cluster-weighted model aims to encompass potential heterogeneity in the distribution of the response variable as well as in the multivariate distribution of the covariates for detecting signals relevant to the underlying latent structure. Of particular interest are potential signals originating from: (i) the linear predictor structures of the regression models and (ii) the covariance structures of the covariates. These two components are modeled using a lasso shrinkage prior for the regression coefficients and a graphical-lasso shrinkage prior for the covariance matrices. A fully Bayesian approach is followed for estimating the number of clusters by treating the number of mixture components as random and implementing a trans-dimensional telescoping sampler. Alternative Bayesian approaches based on overfitting mixture models or using information criteria to select the number of components are also considered. The proposed method is compared against EM type implementation, mixtures of regressions, and mixtures of experts. The method is illustrated using a set of simulation studies and a biomedical dataset.

C0367: **Mixture-based clustering for ordinal responses**

Presenter: **Marta Nai Ruscone**, Università degli Studi di Genova, Italy

Existing methods can perform likelihood-based clustering on a multivariate data matrix of ordinal data, using finite mixtures to cluster the rows (observations) of the matrix. These models can incorporate the main effects of individual rows and columns, as well as cluster effects, to model the matrix of responses. However, many real-world applications also include available covariates, which can provide insights into the main characteristics of the clusters. The mixture-based models are extended to include covariates directly, to allow the clustering structures to be determined both by the individuals' similar patterns of responses and the effects of the covariates on the individuals' responses. The focus is on clustering the rows of the data matrix, using the proportional odds cumulative logit model for ordinal data. The models are fit using the expectation-maximization algorithm, and performance is assessed through a comprehensive simulation study. An application of the models is also illustrated.

C0369: **Mixture of matrix-variate normals with mean restrictions**

Presenter: **Marco Berrettini**, University of Bologna, Italy

Co-authors: Giuliano Galimberti, Cinzia Viroli

Novel strategies are introduced for modeling mean structures in matrix-variate normal distributions, aiming to overcome the common issue of over-parameterization in model-based clustering. The proposed approach relies on parsimonious parameterizations, based on additive components, interaction effects, and low-degree polynomial terms, to effectively capture the intricate multivariate relationships encoded in the mean structure of the data. Particular attention is given to identifiability, and explicit expressions for maximum likelihood estimation under the proposed constraints are derived. To extend these ideas within a clustering framework, finite mixtures of mean-restricted matrix normal distributions are considered, combined with structured covariance matrices that preserve flexibility while reducing the risk of overfitting. An expectation-maximization (EM) algorithm is developed to perform parameter estimation efficiently and reliably. The effectiveness of the proposed methodology is demonstrated through an illustrative application to climate data, where it shows significant improvements over more conventional approaches in detecting meaningful patterns and group structures in high-dimensional, matrix-valued observations.

C0891: **Federated variational inference for Bayesian mixture models**

Presenter: **Paul Kirk**, University of Cambridge, United Kingdom

Co-authors: Jackie Rao

A one-shot, unsupervised federated learning approach is presented for Bayesian model-based clustering of large-scale binary and categorical datasets, motivated by the need to identify patient clusters in privacy-sensitive electronic health record (EHR) data. A principled divide-and-

conquer inference procedure is introduced, using variational inference with local merge and delete moves within batches of the data in parallel, followed by global merge moves across batches to find global clustering structures. It is shown that these merge moves require only summaries of the data in each batch, enabling federated learning across local nodes without requiring the full dataset to be shared. Empirical results on simulated and benchmark datasets demonstrate that the method performs well relative to comparator clustering algorithms. The practical utility of the method is validated by applying it to a large-scale British primary care EHR dataset to identify clusters of individuals with common patterns of co-occurring conditions (multimorbidity).

C0984: Repulsive mixtures via the sparsity-inducing partition prior

Presenter: **Alexander Mozdzen**, A*STAR, Singapore

Co-authors: Gregor Kastner, Andrea Cremaschi, Maria De Iorio, Timothy Wertz

A novel prior distribution is introduced for modelling the weights in mixture models based on a generalization of the Dirichlet distribution, the Selberg Dirichlet distribution. The prior contains a repulsive term, which naturally penalises values that lie close to each other on the simplex, thus encouraging few dominating clusters. The repulsive behaviour induces additional sparsity in the number of components. This construction is referred to as a sparsity-inducing partition (SIP) prior. By highlighting differences with the conventional Dirichlet distribution, relevant properties of the SIP prior are presented, and their implications are demonstrated across a variety of mixture models, including finite mixtures with a fixed or random number of components as well as repulsive mixtures. An efficient posterior sampling algorithm is proposed, and the model is validated through an extensive simulation study as well as an application to a biomedical dataset describing children's Body Mass Index and eating behavior.

CO332 Room BCB 407 ADVANCES IN MODERN CAUSAL INFERENCE

Chair: Mireille Schnitzer

C1397: Using a generalization of the average treatment effect on the treated to better understand opioid prescription risks

Presenter: **Kara Rudolph**, Columbia University, United States

Co-authors: Ivan Diaz, Shodai Inose, Herbert Susmann, Nicholas Williams, Katherine Hoffman, Allison Perry

To reduce risks associated with using prescription opioids, the US (CDC) published opioid prescribing guidelines. However, these guidelines 1) do not consider dose strength and prescription duration as a joint exposure and 2) are written as applying to all persons. However, perhaps the majority of opioid prescribing poses little risk and, therefore, is not in need of intervention. A large cohort of opioid-naïve musculoskeletal pain patients on Medicaid is considered, and effects of modest reductions in opioid dose and duration prescribing practices (considered as a joint exposure) on risk of developing opioid use disorder are estimated, among patients with different prescribing levels. This causal effect is a generalization of the average treatment effect on the treated (ATT); it is the effect of a modified treatment policy among subsets for whom the policy would be relevant, based on their treatment status. A novel targeted minimum loss-based estimator is used with cross-fitting. It is found that universal reductions may have little effect on reducing opioid-related harms, and, plausibly, may even be counterproductive to the extent that their universal application results in uncontrolled pain for patients. In contrast, applying opioid prescribing guidelines to target patients with high-dose and/or long-duration prescriptions would be expected to yield much larger benefits.

C1458: Causal optimal transport of treatment effect to a target population with limited individual-level data

Presenter: **Tat Thang Vo**, University Paris Est Creteil, France

Co-authors: Antoine Chambaz

The transportability of empirical findings to new environments, settings, or populations is essential in most scientific investigations. One practical challenge of standard methods for transportability, however, is that they require the individual-level data on outcome, treatments, and case-mix characteristics to be fully accessible in the source study, along with individual-level data on case-mix characteristics for a random sample from the target population. In practice, data sharing is often subject to administrative barriers and privacy concerns, e.g., a pharmaceutical company may have access to individual-level data from its own study but only aggregate-level data for the target population, such as from a competitor's trial. In such a scenario, state-of-the-art methods generally rely on parametric G-computation or inverse weighting to adjust for the case-mix difference between the study and the target population. Unsurprisingly, the resulting effect estimates from these approaches can be severely biased if the modeling assumption imposed on the nuisance outcome and/or weight model is violated. Novel methods are developed for transportability that fully circumvent the need for strong parametric assumptions when there is restricted access to individual-level data, using computational optimal transport. The new methods allow the use of flexible data-driven methods to estimate nuisance parameters and rely on semi-parametric theory for valid asymptotic inference.

C1462: Causal vaccine effects in the naturally infected

Presenter: **David Benkeser**, Emory University, United States

Establishing the long-term effects of interventions aimed at preventing intermediate outcomes poses challenges. For example, vaccines designed to prevent diarrhea caused by Shigella bacteria in children may also positively impact long-term growth, as Shigella-induced diarrhea is a known cause of growth faltering. However, given the relatively low frequency of Shigella-related diarrhea, the vaccine's marginal causal effect on growth may be too small to detect in a randomized controlled trial. Nevertheless, policymakers are interested in demonstrating the broader benefits of vaccination on growth outcomes. To address this challenge, alternative estimands that enjoy improved power for detecting effects on long-term outcomes in realistic trial settings are proposed. Both principal stratification and interventional causal frameworks are used, and both approaches are demonstrated to yield the same identifying functional under different assumptions. Notably, the principal stratification approach relies on cross-world independence assumptions, whereas the interventional estimand does not. Nonparametric efficient, and doubly robust estimators for these estimands are further derived, leveraging machine learning techniques for nuisance parameter estimation. Through realistic simulations, these estimators are shown to provide clinically meaningful inferences even within the constraints of practical Shigella vaccine trials.

C1468: Mixed-effects and latent space approaches for causal inference under network interference

Presenter: **Vanessa McNealis**, University of Glasgow, United Kingdom

Co-authors: Erica Moodie, Nema Dean

Causal analyses have grown popular, yet most methods rely on a no-interference assumption, in which an individual's outcome is unaffected by the exposures of others. However, some exposures exhibit spillover, affecting both the direct recipient and their social contacts, as in a sexual health promotion intervention delivered in schools. Estimating causal spillover is complicated by homophily, where individuals connect with others sharing similar characteristics, and by unmeasured cluster-level factors. The purpose is to discuss recent advances addressing unmeasured homophily, contextual confounding, and realistic evaluation. First, a Bayesian joint mixed-effects framework for clustered network data is highlighted, that combines outcome and exposure models with direct standardization to enable causal estimation under cluster-level confounding. Second, latent space approaches to construct estimators of direct and spillover effects are described, accounting for uncertainty in latent positions in the social space and allowing flexible outcome models. Finally, developed plasmode simulations are presented to evaluate causal inference methods on realistic social networks. Together, these approaches provide tools for causal inference under interference, accounting for homophily and unmeasured confounding.

C1503: Structural nested recurrent event model for estimating the effects of time-varying exposure

Presenter: **Ashkan Ertefaie**, University of Pennsylvania, United States

Co-authors: Daniel Mork, Francesca Dominici, Robert Strawderman

Assessing the causal effect of time-varying exposures on recurrent event processes is challenging in the presence of a terminating event. The

objective is to estimate both the short-term and delayed marginal causal effects of exposures on recurrent events while accounting for bias introduced by a potentially correlated terminal event. Existing estimators based on marginal structural models and proportional rate models are unsuitable for estimating delayed marginal causal effects for many reasons, and furthermore, they do not account for competing risks associated with a terminating event. To address these limitations, we propose a class of semiparametric structural nested recurrent event models and two estimators of short-term and delayed marginal causal effects of exposures. We establish the asymptotic linearity of these two estimators under regularity conditions through the novel use of modern empirical process and semiparametric efficiency theory. We examine the performance of these estimators via simulation and provide an R package *scure* to apply our methods in real data scenarios. Finally, we present the utility of our methods in the context of a large epidemiological study of 299,661 Medicare beneficiaries, where we estimate the effects of fine particulate matter air pollution on recurrent hospitalizations for cardiovascular disease.

CO342 Room BCB 409 SPATIAL STATISTICS, IMAGE ANALYSIS AND DIRECTIONAL DATA
Chair: Ranjan Maitra
C0946: Matrix-free conditional simulations of Gaussian random fields

Presenter: **Somak Dutta**, Iowa State University, United States

Co-authors: Debashis Mondal

In spatial analysis, conditional simulation of spatial variables at unobserved locations given the data at the observed location facilitates various statistical inferences but suffers from computational scalability when the sample size is large. The aim is to develop a matrix-free method for conditional simulation based on novel mathematical decompositions of the inverse-covariance matrix. The method applies to a broad class of spatial models, including the Gaussian Markov random fields, fractional Gaussian fields, and the Matern models. A practical application is described to mapping groundwater arsenic exceedance regions.

C0947: Approximations in the anisotropic Ising and 3-color Potts models for use in scene analysis

Presenter: **Alejandro Murua**, University of Montreal, Canada

Co-authors: Ranjan Maitra

The Potts distribution is useful in many applications involving statistical modeling and inference. However, its normalization constant and its moments are intractable. This drawback has prevented a wider use of this model. The recent work is extended to fast numerical approximations to the homogeneous Ising model. Analytical approximations are supplied for the anisotropic Ising model and the 3-color Potts model. Simulation studies indicate good performance compared to Markov chain Monte Carlo methods and in a tiny fraction of the time. The methodology is illustrated with real applications and hypothesis testing.

C1032: An adaptive modal EM with application to image segmentation

Presenter: **Giovanna Menardi**, University of Padova, Italy

Image segmentation is the automated process of identifying distinct regions within a digital image for purposes as content retrieval, object detection, or pattern recognition. Digital images are composed of a fixed number of pixels, discrete elements encoding quantized values representing color or intensity. Segmentation involves assigning a label to each pixel so that those sharing similar properties, such as color, intensity, or texture, are grouped together. This goal naturally aligns with cluster analysis, which has become a central tool in image segmentation. Among clustering techniques, nonparametric methods are particularly well-suited for segmentation tasks, as they can identify segments of arbitrary shape and do not require the number of segments to be specified in advance. One of the most widely used nonparametric approaches is the mean-shift algorithm, a gradient ascent procedure based on kernel density estimation, where segments are associated with the modes in the color intensity and possibly spatial distribution. Kernel density estimators can perform poorly when modal regions vary in shape and size. Motivated by this limitation, an alternative nonparametric method is developed for adaptive density estimation and simultaneous group identification, which relies on a two-level EM-style algorithm: One layer for parameter estimation and one for mode detection. The aim is to investigate its application to the problem of segmentation in greyscale images.

C0878: Bridging geometry and statistics: PCA for directional data

Presenter: **Anahita Nodehi**, University of Bristol, United Kingdom

Co-authors: Meisam Moghimbeygi, Christophe Ley

High-dimensional data often present significant challenges in statistical analysis, including difficulties in visualization, increased computational complexity, and a higher probability of overfitting or underfitting. These issues are further compounded by the curse of dimensionality, which states that the number of observations required for accurate modeling grows exponentially as the number of dimensions increases. To address these challenges, dimension reduction techniques are commonly employed. Principal component analysis (PCA) is one of the most widely used dimension reduction techniques and has been extensively studied within classical linear (Euclidean) spaces. However, in many applied fields such as biology, bioinformatics, astronomy, and geology, data often lie in non-Euclidean spaces, specifically on Riemannian manifolds such as the unit circle, sphere, or torus. In these contexts, data are referred to as manifold-valued or directional data. When working with directional data, the linear assumptions underlying PCA may limit its effectiveness for accurate dimension reduction. The purpose is to review and investigate the methodological developments of PCA for directional data and explore their practical applications.

C1033: A semiparametric quasi-likelihood regression for circular responses

Presenter: **Anna Gottard**, University of Firenze, Italy

Co-authors: Andrea Meilan-Vila, Agnese Panzera

Flexible and interpretable methods for circular outcomes are in high demand. A semiparametric regression framework is proposed for a circular response that accommodates both linear and circular predictors in its parametric and nonparametric components. Instead of assuming a specific error distribution, a circular quasi-likelihood is employed. The proposed semiparametric regression method bridges the gap among fully parametric circular regression, fully nonparametric kernel methods, and spline-based GAMs with cyclic smoothers, offering both flexibility and interpretability. The proposed regression method can be adopted in several fields, such as meteorology, for studying wind and wave direction, ecology, for animal movement and migration direction, and any context involving angular or periodic data. The asymptotic behavior of the estimators is established, and a backfitting algorithm is outlined for estimation.

CC350 Room BCB 312 TIME SERIES INFERENCE AND ESTIMATION
Chair: Joachim Schnurbus
C0246: Bias-reduction in state-space model estimation

Presenter: **Magda Monteiro**, University of Aveiro, Portugal

Co-authors: Marco Costa

State-space models (SSM) provide a versatile framework for representing dynamic systems, where latent states evolve through a stochastic process and observed data are connected via a linear measurement equation. A method is proposed for estimating the parameters of SSM, based on a double-iterated generalized method of moments (GMM) procedure. The proposed approach is designed to reduce estimation bias, particularly in contexts involving small sample sizes or high levels of variability. Simulation results demonstrate that this method outperforms traditional estimators, such as maximum likelihood, in terms of accuracy and robustness.

C1282: The joint asymptotic distribution of permutation entropy and complexity*Presenter:* **Angelika Silbernagel**, Helmut Schmidt University, Germany*Co-authors:* Christian Weiss

Since they were first introduced, ordinal patterns have gained popularity as a method for data analysis. As the term suggests, ordinal patterns capture the ordinal structure of the underlying data. They have many desirable properties, like invariance under monotone transformations, robustness with respect to small noise, and extremely fast calculation. Later, the permutation entropy-complexity pair, which is based on ordinal patterns, emerged as a popular tool for summarizing the time-series dynamics. To best of knowledge, one of the main gaps of this tool so far is the lack of a theoretical foundation for its estimation uncertainty. Although deriving the exact (joint) sampling distribution of entropy and complexity is hardly possible, a promising approach is to make use of its asymptotic properties. Therefore, the asymptotic distribution of the entropy-complexity pair is deduced under mild dependence conditions, making the necessary distinction between a uniform and a non-uniform ordinal pattern distribution. In that way, two different limit theorems are obtained. Finally, the theory is complemented by proposing methods for visualizing the estimation uncertainty.

C1261: From spectral decomposition to projection: Rethinking common components in dynamic factor models*Presenter:* **Jan Bruha**, Czech National Bank, Czechia

The purpose is to revisit the estimation of the common component in dynamic factor models. Standard frequency-domain approaches reconstruct the common part via the inverse Fourier transform of spectral eigenvectors, which yields a two-sided filter. One-sided alternatives reduce to filtering conditional on observations at one particular date point, which, for some processes, may result in not efficiently capturing the lead-lag structure among series. A projection-based method is proposed that estimates the common component directly, using the spectral density to recover its autocorrelation function and cross-correlations with observables. This allows filtering based on any available data, even at the sample boundaries or in the presence of missing values due to publication lags. In addition, it is shown how this projection approach permits economically meaningful restrictions on filter weights, with an application to underlying inflation indicators.

C1312: Self-normalization for CUSUM-based change detection in locally stationary time series*Presenter:* **Florian Heinrichs**, FH Aachen - University of Applied Sciences, Germany

A novel self-normalization procedure for CUSUM-based change detection in the mean of a locally stationary time series is introduced. Classical self-normalization relies on the factorization of a constant long-run variance and a stochastic factor. In this case, the CUSUM statistic can be divided by another statistic proportional to the long-run variance, so that the latter cancels. Thereby, a tedious estimation of the long-run variance can be avoided. Under local stationarity, the partial sum process converges to $\int_0^t \sigma(x) dB_x$ and no such factorization is possible. To overcome this obstacle, a self-normalized test statistic is constructed from a carefully designed bivariate partial-sum process. Weak convergence of the process is proven, and it is shown that the resulting self-normalized test attains asymptotic level α under the null hypothesis of no change, while being consistent against a broad class of alternatives. Extensive simulations demonstrate better finite-sample properties compared to existing methods. Applications to real data illustrate the method's practical effectiveness.

C1154: Nonstationarity extended Whittle estimation of cyclical time series*Presenter:* **Edward Hill**, Queen Mary University of London, United Kingdom

The Gegenbauer ARMA (GARMA) process is a fractional differencing model fitted to time series containing a strongly dependent cycle. In many real-world applications, the length of the cycle is unknown, and it is not clear whether the data is stationary. A definition of a GARMA process that also covers nonstationarity is provided. The model is characterised by the cyclicity parameter ω , memory parameter d , and short memory AR and MA parameters. A two-stage procedure to estimate all the model parameters is provided. First, the cyclicity parameter ω is estimated by maximising the periodogram over the Fourier frequencies. The n -consistency and limit distribution of the estimator of ω are obtained. Next, the memory parameter d and AR and MA parameters are estimated using a modified Whittle likelihood that uses the cycle frequency from the first stage estimator. The second stage estimates are shown to be consistent and asymptotically normal for both stationary and nonstationary GARMA processes. The complete estimation procedure allows for the construction of confidence intervals for all the GARMA model parameters, and Monte Carlo simulations confirm its good finite sample performance.

CC357 Room BCB 313 SURVIVAL ANALYSIS**Chair: Angeles Carnero****C1110: Dynamic prediction for the joint model of a longitudinal biomarker and interval-censored failure time data***Presenter:* **Yang-Jin Kim**, Sookmyung Women University, Korea, South

Joint modeling approaches have been widely used to reflect the effect of time-varying biomarkers on failure time. The aim is to propose a joint model for interval-censored failure time data in the presence of longitudinal biomarkers. To assess predictive accuracy, dynamic measures are introduced, including the time-dependent area under the ROC curve (AUC) and the Brier score. The inference procedure is based on the EM algorithm that accounts for the latent failure time and subject-specific random effects. A dynamic marker is defined as the conditional survival probability of being alive at time $t > s$, given survival up to time s . Various simulation scenarios are considered to validate the proposed joint model and compare its performance to the landmarking approach. Finally, the method is illustrated using the well-known Paquid dataset, which includes interval-censored dementia onset times and two longitudinal cognitive scores.

C1298: Penalized piecewise exponential distributional regression model for survival analysis*Presenter:* **Jack Moore**, University of Limerick, Ireland*Co-authors:* Shirin Moghaddam, Kevin Burke

The field of survival analysis is concerned with the modelling of time-to-event data, a form of data that arises in various application areas. A key application area of survival analysis is medical research, where interest lies in survival times of patients. Traditional parametric modelling approaches rely on distributions such as the Weibull, which can impose strong assumptions on the data at hand. In contrast, the piecewise exponential model offers a much more general parametric modelling approach: It has the capability of approximating any survival distribution, without prior knowledge of the underlying distribution. This is achieved by partitioning the time scale into intervals, within which the hazard rate is constant. Despite the versatility of the piecewise exponential model, it has historically been underutilized within the literature. This can perhaps be explained by the popularity of the Cox model. However, in recent years, there has been renewed interest in the piecewise exponential model, with various developments aimed at enhancing its viability. The purpose is to introduce the piecewise exponential model and present extensions aimed at improving this model. Specifically, a distributional regression structure is used. Thus, the framework enables flexible modelling of both the underlying baseline hazard and the nature of covariate effects, where the intervals/ pieces and covariates of the model are selected automatically via an adaptive lasso penalization.

C1374: Monitoring time to event in medical registry data using CUSUMs based on excess hazard models*Presenter:* **Jan Terje Kvaloy**, University of Stavanger, Norway*Co-authors:* Jimmy Huy Tran, Hartwig Korner

In health registries, patients are routinely registered at diagnosis, and outcome data like survival times are added later. Based on such data, an aspect of interest could be to monitor whether the distribution of the time to an outcome of interest changes over time, for instance, if the survival time distribution of cancer patients changes over time, while adjusting for known risk factors. A challenge in monitoring survival times based on such registry data is that the cause of death is often not registered. To quantify the burden of disease in such cases, excess hazard methods, where the

total hazard is modelled as the population hazard plus the excess hazard due to the disease, can be used. The aim is to propose a CUSUM procedure for monitoring for changes in the time to event distribution in such cases. The procedure is based on a survival loglikelihood ratio which extends previously suggested methods for monitoring of time to event to the excess hazard setting. The procedure considers changes in the population risk over time, as well as changes in the excess hazard explained by observed covariates. Properties, challenges, and an application to cancer registry data are presented.

C1432: Joint modeling of zero-inflated longitudinal and survival data with a cure fraction: An application to AIDS data

Presenter: **Taban Baghfalaki**, The University of Manchester, United Kingdom

Co-authors: Mojtaba Ganjali

The purpose is to develop a Bayesian joint modeling framework for analyzing zero-inflated longitudinal count data and survival outcomes, with explicit incorporation of a cure fraction to account for individuals who will never experience the event of interest. The longitudinal trajectory is modeled using a flexible mixed-effects Hurdle specification to address excess zeros and overdispersion that commonly arise in biomedical count data. The survival process is represented through a Cox proportional hazards mixture cure model, enabling separation of cured from susceptible subjects. To capture the association between the two processes, the survival model incorporates a linear combination of current longitudinal values as time-dependent predictors. Bayesian inference is carried out using Hamiltonian Monte Carlo, which ensures efficient posterior sampling and reliable parameter estimation in complex settings. The framework supports dynamic prediction, enabling individualized risk assessment and personalized clinical decision-making. Extensive simulation studies are conducted to evaluate estimation accuracy, predictive performance, and robustness of the proposed model. Finally, the methodology is applied to a real AIDS cohort, illustrating its capacity to integrate longitudinal biomarkers with survival information for improved prediction of patient outcomes. Results underscore the clinical utility of joint modeling in advancing personalized medicine.

C1478: On the choice of censoring distribution in the simulation of survival data

Presenter: **Takis Besbeas**, Athens University of Economics and Business, Greece

Co-authors: Konstantinos Pateras

Simulation of censored time-to-event data is important but presents its own challenges. First, the simulation of right-censoring, which is the most common censoring mechanism in medical research, requires, in general, simulation from two processes, one describing time-to-event and the other time-to-censoring. Second, it will often be desirable to control the censoring proportion of the simulated data. And third, there are various characteristics that make survival data more complex than other types of data. We consider the choice of censoring distribution in the simulation of censored data from a known survival model. We show that there are often practical advantages in selecting the censoring distribution to be in the same probability distribution family as the lifetime distribution, and we illustrate using Monte Carlo simulation that the choice of distribution with bounded support, such as the uniform, which is typically used, may have disadvantages compared to distributions with unbounded support, such as the exponential. We also evaluate the performance of more flexible censoring distributions, such as the Weibull and Gamma, under a range of values for their additional parameter. We further illustrate the effect of the censoring distribution on a parametric bootstrap analysis of a real survival data set on shunts in infants with heart disease.

CC367 Room BCB 402 BIostatistics

Chair: Shanshan Ren

C0874: Spatially regularized Gaussian mixtures for clustering spatial transcriptomic data

Presenter: **Andrea Sottosanti**, University of Padova, Italy

Co-authors: Davide Risso, Francesco Denti

Spatial transcriptomics measures the expression of thousands of genes in a tissue sample while preserving its spatial structure. These technologies have enabled the investigation of the spatial variation of gene expressions and their impact on specific biological processes. Identifying genes with similar expression profiles is of utmost importance, thus motivating the development of flexible methods leveraging spatial data structure to cluster genes. A modeling framework is proposed for clustering observations measured over numerous spatial locations via Gaussian processes. Rather than specifying their covariance kernels as a function of the spatial structure, it is used to inform a generalized Cholesky decomposition of their precision matrices. This approach prevents issues with kernel misspecification and facilitates the estimation of a non-stationarity spatial covariance structure. Applied to spatial transcriptomic data, the model identifies gene clusters with distinctive spatial correlation patterns across tissue areas comprising different cell types, like tumor and stromal areas.

C0893: Bi-clustering RNA-sequencing data: A model-based approach using a MPLN distribution

Presenter: **Caitlin Kral**, Carleton University, Canada

Co-authors: Sanjeena Dang, Ryan Browne, Evan Chance

Bi-clustering is a technique that simultaneously clusters observations and features (i.e., variables) in a dataset. This technique is used in bioinformatics to simultaneously identify clusters of disease and non-diseased patients and the network of genes with distinct correlation patterns based on their gene expression values. While several Gaussian mixture models-based biclustering approaches currently exist in the literature for continuous data, approaches to handle discrete data have not been well researched. Extending bi-clustering approaches to discrete data is imperative, as such data is commonly found within real-world applications such as bioinformatics. Recently, multivariate Poisson-lognormal (MLPN) models have emerged as an efficient model for modelling multivariate count data. It arises from a hierarchical Poisson structure, which allows for over-dispersion and correlation (both positive and negative). A MPLN model-based bi-clustering approach that utilizes a block-diagonal covariance structure is proposed. The clustering performance of the proposed model for clustering both observations and features using simulated and real-world data is demonstrated.

C1264: Signal detection in adverse drug reactions under fluctuating reporting rates

Presenter: **Tatsuhiko Anzai**, Institute of Science Tokyo, Japan

Co-authors: Kunihiko Takahashi

Spontaneous adverse event reporting systems, such as the Japanese adverse drug event report database (JADER) and the FDA adverse event reporting system (FAERS), are central resources for pharmacovigilance. A standard approach for signal detection is the reporting odds ratio (ROR), based on disproportionality analysis in a two-by-two contingency table. Conventional methods assume that reporting rates remain stable across drugs and events, but this assumption is frequently violated. During the COVID-19 pandemic, substantial shifts in reporting behavior were observed, leading to spurious increases in disproportionality measures. A statistical framework is proposed that incorporates reporting rate variations into the signal detection process. The method extends the contingency table with four fluctuation parameters representing deviations in cell-specific reporting ratios. These parameters are estimated by minimizing a divergence between observed and predicted counts, with predictions derived from a regression model assuming stable reporting rates. A two-layer error structure accounts for uncertainty in both predicted frequencies and observed data. The proposed method is applied to real pharmacovigilance data, focusing on psychotropic drugs where reporting rates are particularly unstable.

C1488: Bayesian inference for cluster-randomized trials with multivariate outcomes with missing data

Presenter: **Guangyu Tong**, Yale University, United States

Cluster-randomized trials (CRTs) with fragile populations often face complex attrition, in which missing outcomes arise from heterogeneous causes: participants may be alive, deceased, or of unknown status, each with distinct missing-data mechanisms. Existing methods address death

truncation but cannot jointly handle dropout unrelated to mortality or unknown survival. We propose a Bayesian framework to estimate survivor average causal effects in CRTs while accounting for multiple missingness types. Our approach models multivariate outcomes, producing posterior estimates that distinguish individual- and cluster-level survivor effects. Simulation studies demonstrate low bias and high coverage across varied scenarios. We illustrate the method using data from a geriatric CRT, focusing on a bivariate continuous outcome, though the framework readily extends to multiple or alternative endpoints (e.g., binary). We offer a general modeling strategy for handling complex missingness in CRTs, broadly applicable in aging and palliative care research.

C1516: **Confirming kappa-cover: A hypothesis test for the overlap of normal distributions based on generalized p-values**

Presenter: Vera Hofer, University of Graz, Austria

Co-authors: Gerhard Goessler, Hans Manner, Walter Goessler

A common task in various industrial applications is to compare quality characteristics of two products in terms of their distributions, DT and DR. Simply comparing their means and variances is often too restrictive, especially when a certain scope of DT is allowed within the range given by DR. At the same time rather strict control of extreme realizations of DT is often demanded, i.e., DR must cover DT such that it ensures the test product being safe regarding the quality characteristic under investigation. This is reflected by the concept of kappa-cover. The problems associated with developing a suitable test for confirming kappa-cover are discussed. Challenges arise from the multiple testing problem of comparing two pairs of quantiles simultaneously, as well as from the formulation of the hypotheses. From a consumer safety perspective, it is preferable to formulate H_0 as the assertion of missing kappa-cover. However, this raises the question of how to handle distributional equality, which formally falls under H_0 (i.e., no kappa-cover) even though DR evidently covers DT in this situation. The proposed two-stage test procedure ensures that the chosen significance level is asymptotically maintained in this multiple testing problem. According to simulation results, the test procedure controls the type I error in the expected manner, i.e., as the sample size increases, the size and power of the test improve.

CC363 Room BCB 406 NONPARAMETRIC METHODS

Chair: Debarghya Ghoshdastidar

C1095: **Nonparametric welfare analysis for discrete choice with multiple compensation timings**

Presenter: Koshi Nishida, University of Tokyo, Japan

The purpose is to consider welfare analysis in a discrete choice framework where income receipts and price payments occur over multiple periods, allowing for unrestricted unobserved heterogeneity in individual preferences. For changes in prices over multiple periods, the definitions of welfare measures are extended, such as compensating variation and equivalent variation, to incorporate income adjustments over time. This framework naturally generalizes previous works by accommodating multiple periods with fully unrestricted heterogeneity in factors such as borrowing constraints, interest rates, and saving behaviors, all of which are assumed to be unobservable to the econometrician. The situations are then clarified, in which the distribution functions of the welfare measures can or cannot be nonparametrically point-identified.

C1232: **Regularized maximum likelihood estimation for the random coefficients model**

Presenter: Fabian Dunker, University of Canterbury, New Zealand

Co-authors: Emil Mendoza, Marco Reale

A popular way to model unobserved heterogeneity in population is the linear random coefficient model $Y_i = \beta_{i,1}X_{i,1} + \beta_{i,2}X_{i,2} + \dots + \beta_{i,d}X_{i,d}$. It is assumed that the observations (\mathbf{X}_i, Y_i) , $i = 1, \dots, n$, are i.i.d. where $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,d})$ is a d -dimensional vector of regressors. The random coefficients $\beta_i = (\beta_{i,1}, \dots, \beta_{i,d})$, $i = 1, \dots, n$ are unobserved i.i.d. realizations of an unknown d -dimensional distribution with density f_β independent of \mathbf{X}_i . The aim is to propose a quasi-maximum likelihood method to estimate the joint density distribution of the random coefficients. This method implicitly involves the inversion of the Radon transformation in order to reconstruct the joint distribution, and hence is an inverse problem. To add stability to the solution, Tikhonov-type regularization methods are applied. Nonparametric estimation for the joint density of β based on kernel methods or Fourier inversion has been proposed in recent years. Most of these methods assume a heavy-tailed design density $f_{\mathbf{X}}$. The convergence of the quasi maximum likelihood method is analyzed without assuming heavy tails for $f_{\mathbf{X}}$, and performance is illustrated by applying the method on simulated and real data.

C1250: **An optimal transportation approach of confidence intervals**

Presenter: Christophe Valvason, University of Geneva, Switzerland

Co-authors: Stefan Sperlich, Eustasio del Barrio

Reliable inferential tools for small samples and complex statistics are crucial in empirical research and official statistics. The aim is to propose a novel approach to constructing confidence intervals based on optimal transport theory. While in the univariate case, Monge's problem reduces to the composition of the CDF and quantile function, small samples yield empirical CDFs with large jumps, making the standard transport map suboptimal. To address this, the optimal transport plan is computed between a reference distribution and a nonparametric estimate of the distribution of interest. Confidence intervals are first defined in the reference space and then transported back to the original problem. Optimal transport theory provides the theoretical foundation for the validity of this method. Simulation studies demonstrate that the approach achieves coverage probabilities closer to the nominal level and often reduces variance compared to both direct and bootstrap confidence interval estimators.

C1439: **Estimation of the invariant measure of a multidimensional diffusion from noisy observations**

Presenter: Gregoire Szymanski, Universita du Luxembourg, DMATH, France

Co-authors: Raphael Maillet

The purpose is to present a new approach for estimating the invariant density of a multidimensional diffusion when dealing with high-frequency observations that are blurred by independent noise. The focus is on the intermediate regime, where observations are collected at discrete times $k\Delta_n$ for $k = 0, \dots, n$, under the conditions $\Delta_n \rightarrow 0$ and $n\Delta_n \rightarrow \infty$. The method relies on a kernel density estimator combined with a pre-averaging technique, which effectively removes noise from the data while preserving the analytical structure of the underlying signal and its asymptotic behavior. How the rate of convergence of the estimator depends on both the anisotropic regularity of the density and the noise intensity is discussed. In particular, conditions on the noise level that allow achieving convergence rates comparable to the noiseless case are described.

C1491: **Nonparametric kernel mixed data sampling models**

Presenter: Deliang Dai, Linnaeus university, Sweden

A novel extension of the Fully Non-Parametric MIDAS model is introduced. Our approach leverages kernel methods to create a purely nonparametric framework, where estimation complexity depends solely on sample size. This allows the model to capture complex nonlinearities in both the temporal lag structure and the functional relationship between mixed-frequency series. A key feature is the application of a single kernel function to all lags of a given regressor. For estimation, we employ a kernel-based trend filtering algorithm that provides local adaptivity at a lower computational cost than standard spline regression. Empirical results on simulated and real-world data, including a case study forecasting urban air pollution from meteorological conditions, demonstrate the superior performance of our method over linear MIDAS.

CC358 Room BCB 408 RISK ANALYSIS

Chair: Genaro Sucarrat

C1125: **Can uncertainty affect extreme events in the oil market? A MIDAS touch to dynamic POT models**

Presenter: Katarzyna Bien-Barkowska, Poznan University of Economics and Business, Poland

A novel econometric framework is introduced for forecasting extreme events in oil markets by incorporating macro-financial uncertainty indicators

into a dynamic peak-over-threshold (POT) model within a mixed-frequency setting. Based on three decades of daily returns from two major global oil benchmarks, WTI and Brent, robust empirical evidence that rising uncertainty is significantly associated with the frequency and magnitude of extreme returns is provided. The Economic Policy Uncertainty index, Equity Market Volatility Tracker: Commodity Markets, and Jurado-Ludvigson-Ng 1-month ahead macroeconomic uncertainty index are utilized to capture different dimensions of uncertainty. These insights are derived through the novel framework, the autoregressive conditional MIDAS-POT (AC-MIDAS-POT), which integrates dynamic specifications for inter-exceedance times and magnitudes of extreme oil returns with the MIDAS components to account for macroeconomic and financial information. The AC-MIDAS-POT model significantly outperforms the VaR forecasts from the GARCH-MIDAS models with Gaussian and Student's t-distributions. This novel modeling framework offers a robust tool for forecasting tail risk in financial markets and provides valuable insights for investors and policymakers to forecast and mitigate the impact of extreme returns.

C1134: The value-at-risk of a large insurance portfolio: Accounting for model risk

Presenter: **Hanieh Amjadian**, Concordia University, Canada

Co-authors: Yang Lu, Yunran Wei

Recently, quantifying model risk on risk measures such as value-at-risk (VaR) for portfolios has been extensively studied in the context of credit risk, as well as in market risk, particularly within a time series framework. However, relatively few frameworks address model risk in a cross-sectional setting. One notable exception is a recent study, which explicitly considers model risk concerning the aggregation of individual risks. The purpose is to adopt their methodology in quantifying model risk when calculating risk measures, such as VaR or ES of a large insurance portfolio, and address their limitations by introducing an additional calibration dataset. Simulation studies are conducted using heavy-tailed data to evaluate the performance of our proposed method. The results show that relying solely on training data can lead to significant underestimation of risk, particularly when models are overfitted or misspecified. By incorporating a calibration dataset and using a bootstrap procedure, the approach produces more accurate and robust VaR estimates, especially in the tails of the distribution. It is also investigated how model complexity, the heaviness of the tails, and the number of bootstrap replicates affect the quality of the estimates, offering practical guidance for risk assessment in insurance portfolios.

C1248: Corporate default prediction under distribution shift

Presenter: **Suguru Yamanaka**, Aoyama Gakuin University, Japan

Distribution shifts, or covariate shifts, arising from changing economic environments pose a significant challenge to the estimation of corporate default prediction models. A default prediction model based on an importance-weighted least-squares probabilistic classifier (IWLSPC) is introduced to address this degradation in model performance. In contrast to the standard least-squares probabilistic classifier (LSPC), IWLSPC adapts to these shifts by incorporating data density ratios as importance weights in model estimation. Empirical validation using real corporate default datasets demonstrates that the default prediction model based on IWLSPC achieves significantly higher predictive accuracy than the conventional, unweighted LSPC. The findings confirm that importance weighting is a robust technique for building reliable credit risk models capable of operating under non-stationary conditions.

C1355: Risk budgeting portfolios under different risk measures: A comparison of alternative computational approaches

Presenter: **Caterina Pastorino**, University of Milano - Bicocca, Italy

Co-authors: Akif Ince, Pierpaolo Uberti, Ilaria Peri

Modern portfolio theory has provided quantitative tools for decades to help investors make better decisions in financial markets. In recent years, a wide range of portfolio construction methods has been developed. Among risk-based approaches, the risk budgeting portfolio has gained popularity. Rather than focusing on the overall portfolio risk, risk budgeting emphasizes the contribution of each asset to the total portfolio risk. It relies on a mathematical framework to decompose the total portfolio risk into individual asset risk contributions. The aim is to investigate risk budgeting portfolios by comparing several computational methods: Traditional optimization techniques, non-linear system and fixed-point algorithm. These approaches are examined in terms of computational efficiency and accuracy. Moreover, the analysis is conducted using different risk measures: Standard deviation, value at risk (VaR), and expected shortfall (ES). The preliminary findings provide valuable information about the performance of the computational methods implemented.

C1474: Disaster impact forecasting framework for multi hazard disaster risk assessment

Presenter: **Mohammad Reza Yeganegi**, International Institute for Applied Systems Analysis (IIASA), Austria

Co-authors: Nadejda Komendantova

Disaster risk management relies largely on the estimated disaster impact under different scenarios, as well as the estimation of the disaster probability itself. The results of such disaster risk estimation models are inputs to disaster risk management strategy building and decision-making. Decision-makers need an estimation of the risk metrics under each scenario to determine the best combination of strategies for managing disaster risk. Various studies provided models for estimating the probabilistic behavior of the disastrous event. However, the vast diversity of risk metrics and risk drivers poses a challenge in forecasting the disaster impact and consequently, the disaster risk. Furthermore, disaster risk drivers (e.g., public trust, social vulnerability, socio-economic variables, climate change) have a dynamic nature, and it is crucial to consider their dynamic structure when estimating the disaster impact. In other words, the impact should be estimated considering the interactions among the risk metrics as well as adapting to the changes in risk drivers. Furthermore, the comprehensive risk assessment relies on the distribution estimation of the disaster impacts. The purpose is to formulate the problem of the disaster impact estimation in relation to the disaster risk assessment and decision-making, and propose a framework for forecasting the disaster impact. The proposed solution extends the commonly used value at risk (VaR) concept for disaster impact forecasting.

Saturday 13.12.2025

13:40 - 15:20

Parallel Session C – CFE-CMStatistics 2025

CI017 Room BCB 307 ECOSta PART B: STATISTICS INVITED SESSION (VIRTUAL)**Chair: Cristian Gatu****C0169: Physics informed statistical learning for spatial and functional data***Presenter:* **Laura Sangalli**, Politecnico di Milano, Italy

Physics-informed statistical learning methods are a new class of nonparametric and semiparametric regression models with roughness penalties. These methods can handle spatial and spatiotemporal data, as well as functional data, observed over multidimensional domains that can have complex shapes, such as non-convex planar domains, curved domains, non-convex volumes, and linear networks. By integrating differential operators, ranging from simple second derivatives and Laplacians to more sophisticated partial differential equations, these models leverage problem-specific knowledge to enhance model accuracy. The use of unstructured mesh discretization enables high modeling flexibility, making it possible to capture highly localized signals, strong anisotropies, and non-stationary patterns. The utility of these methods are illustrated through applications to complex data analysis problems from life and environmental sciences.

C0244: Machine learning and statistical methods for high-throughput experimental data*Presenter:* **Tatyana Krivobokova**, University of Vienna, Austria

Machine learning (ML) and artificial intelligence (AI) techniques are transforming the way chemical reactions are studied today. Valuable datasets from high-throughput experimentation (HTE) are increasingly being generated to better understand reaction conditions that are crucial for outcomes such as yields and selectivities. However, it is often overlooked that data from such designed experiments possess a very specific structure, which can be captured by appropriate statistical models. Ignoring these underlying data structures when applying ML/AI algorithms can result in completely misleading conclusions. In contrast, leveraging knowledge about the data-generating process together with suitable estimation approaches yields reliable, interpretable, and comprehensive insights into the chemical reaction mechanisms. A particularly complex dataset is available for the Buchwald-Hartwig amination. Using this dataset, an appropriate statistical model for HTE-generated chemical data is introduced, and a suitable parameter estimation algorithm is developed. Based on the estimated model, new insights into the Buchwald-Hartwig amination are thoroughly discussed. The approach is directly applicable to a wide range of HTE-generated data for chemical reactions and beyond.

C0899: Model-free bootstrap and conformal prediction in regression*Presenter:* **Dimitris Politis**, University of California, San Diego, United States*Co-authors:* Yiren Wang

Predictive inference under a general regression setting is gaining more interest in the big-data era. In terms of going beyond point prediction to develop prediction intervals, two main threads of development are conformal prediction and Model-free prediction. The two are contrasted with an emphasis on conditional inference.

CO186 Room BCB G07 HiTEC: SOME NEW CHALLENGES IN FUNCTIONAL STATISTICS**Chair: Enea Bongiorno****C0661: Multi-object regression: A linear framework using partial least squares***Presenter:* **Robert Cantwell**, University of Cambridge, United Kingdom*Co-authors:* John Aston

Modern statistical problems regularly involve data collected with more structure than simple scalars, for example, functional or image data. Accompanying these more exotic data types has been the generalization of statistical methods to analyze data that take values in almost arbitrary Hilbert spaces. However, to date, nearly all regression approaches treat a single, often specific, data type within a single Hilbert space. A linear framework for statistical analysis of data objects naturally represented in a variety of Hilbert spaces, using latent space projections, is considered. In particular, the method of partial least squares are generalized, and significance testing of different data objects is discussed to conduct model selection, a challenge often raised in applications where complicated data objects can be costly to collect and practitioners use statistical modeling to understand the relative importance of each object.

C1078: Goodness-of-fit and distribution tests for functional data*Presenter:* **Daniel Hlubinka**, Univerzita Karlova, Czech Republic

The aim is to present a method for constructing tests using empirical characteristic functionals to test the distribution of functional random variables. In particular, a goodness-of-fit test and a test of symmetry and time reversibility for continuous stochastic processes are presented. Characteristic functionals can be used to construct distance-based tests for distribution or distributional properties of functional data. This approach is particularly useful because the characteristic functional can be consistently estimated, and test statistics are derived as if the functions were fully observable. Once constructed, the test statistics are adjusted because the underlying functions are only observed at a discrete grid. Typically, the Cramer-von Mises test statistic is employed, which is based on the integrated distance between the empirical characteristic functional and its counterpart under the null hypothesis. As the exact distribution of the test statistic under the null hypothesis is unknown and the asymptotic distribution is usually very complicated, some resampling or subsampling methods are used to derive the critical values for the test. The choice of probability distribution for the Cramer-von Mises statistic is also discussed, and the performance of the tests is demonstrated under the null hypothesis and various alternatives.

C1087: Detecting covariance shifts in multichannel profiles*Presenter:* **Davide Forcina**, University of Naples Federico II, Italy*Co-authors:* Christian Capezza, Antonio Lepore, Biagio Palumbo

Modern industrial systems generate multichannel profile data continuously, requiring effective real-time monitoring and fault diagnosis. While many existing methods prioritize detecting shifts in the process mean, changes in the covariance structure are just as important, as they reflect the dynamic interdependencies among multiple variables. A functional graphical modeling framework is introduced to represent conditional dependencies in multichannel profile data, addressing challenges posed by high dimensionality and sparsity. The method leverages penalized likelihood ratio tests with adaptive penalty terms to detect a wide range of covariance structure changes. To enhance interpretability, a diagnostic procedure based on change-point detection is used to pinpoint the specific relationships that have changed. Simulation studies and a case study on multichannel temperature profile monitoring demonstrate the superior performance of the proposed approach compared to existing methods.

C1098: Kernel ridge regression for spherical responses*Presenter:* **Almond Stoecker**, Ecole polytechnique federale de Lausanne, Switzerland*Co-authors:* Beatrice Matteo, Shahin Tavakoli

The aim is to propose a novel nonlinear regression framework for responses taking values on a hypersphere. Rather than performing tangent space regression, where all the sphere responses are lifted to a single tangent space on which the regression is performed, we estimate conditional Frechet means by minimizing squared distances on the nonlinear manifold. Yet, the tangent space serves as a linear predictor space where the regression function takes values. The framework integrates Riemannian geometry techniques with functional data analysis by modelling the regression function using methods from vector-valued reproducing kernel Hilbert space theory. This formulation enables the reduction of the infinite-dimensional estimation problem to a finite-dimensional one via a representer theorem and leads to an estimation algorithm by means of

Riemannian gradient descent. Explicit checkable conditions on the data that ensure the existence and uniqueness of the minimizing estimator are given.

CO236 Room BCB G08 RECENT CONTRIBUTIONS TO RISK ANALYSIS
Chair: Angeles Carnero
C0458: Functional modeling of electricity returns with dynamic volatility

Presenter: **Eric Costa Andreu**, Universidad de Alicante, Spain

Co-authors: Angeles Carnero, Pedro Galeano

A functional time series model is used to capture the dynamics of daily return curves of electricity prices. While electricity prices are typically recorded hourly as discrete values, representing them as continuous functions could better capture the underlying structure of price fluctuations. Using data from the Spanish electricity market, functional representations of daily returns are constructed, their temporal stability is assessed, and functional principal component analysis (FPCA) is applied. Components accounting for a large amount of the total variability are retained, and their scores are individually modeled using a seasonal ARMA-GJR-GARCH process estimated through maximum likelihood. Predicted scores enable the reconstruction of the return curves. Moreover, the proposed functional model facilitates the estimation and forecasting of return volatility, providing a natural measure of market risk. Lastly, the statistical properties of the model's parameter estimators are examined by analyzing their empirical distributions using both parametric and residual-based bootstrap techniques.

C0569: Asymmetric effects on asymmetry: The resilience of ESG indices

Presenter: **Javier Perote**, University of Salamanca, Spain

Co-authors: Ines Jimenez, Andres Mora Valencia

A recently developed multivariate volatility model is applied based on multivariate Gram-Charlier (GC) expansions to investigate spillover effects between traditional and environmental, social, and governance (ESG) based indices. The model extends the GC distribution by incorporating the crossed products of Hermite polynomials among different moments and assets in a multivariate expansion. The weights of these new terms are argued to convey spillover effects and have been proved to be significant pieces of information for tail dependence measurement. The empirical research analyzes bivariate portfolios of S&P500 and different ESG-based exchange trade funds (ETFs). The evidence shows that ESG ETFs are not immune to clusters of high volatility transmission (cross-kurtosis spillover), but governance and socially concerned assets seem to be more resilient to transitory shocks (cross-skewness spillover). Interestingly, clean energy and environmentally focused ETFs present significant cross-skewness linkages with traditional market indices.

C0699: Expected shortfall estimation with stationary vine copulas

Presenter: **Juan Mora**, Universidad de Alicante, Spain

Co-authors: Roberto Fuentes-Martinez

The performance of stationary vine copula models (S-vines) is assessed for estimating the expected shortfall (ES) of returns from a portfolio of financial assets. To this end, an estimation procedure is proposed based on the Monte Carlo method using S-vines, for estimating the k-period-ahead ES of a financial portfolio at a given time period t , employing a rolling window approach. Notably, by means of a simulation study, evidence is found that, under some dependence scenarios, the S-vine ES estimates outperform those of models commonly used in financial time series modeling, while exhibiting relatively similar performance under other data-generating processes. Finally, through an empirical application, it is shown that using the S-vine ES estimates in the context of portfolio optimization can lead to better strategies than those derived from models that consider serial dependence and cross-sectional dependence individually, particularly when working with large portfolios.

C0731: New robust conditional quantile-based volatility, skewness and kurtosis

Presenter: **Trino Niguez**, University of Westminster, United Kingdom

Co-authors: Laura Garcia-Jorcano, Angel Leon

The aim is to introduce a novel and computationally efficient method to estimate time-varying skewness and kurtosis of daily stock returns using quantile-based techniques. The proposed two-step approach first estimates return quantiles using the CAViaR model, then fits a cubic polynomial to derive conditional skewness and kurtosis via non-linear transformations of the estimated coefficients. Classical quantile-based measures of unconditional skewness and kurtosis are revisited, and their limitations are highlighted relative to moment-based definitions. Recent contributions propose alternative quantile-based measures aligned with moment theory. Building on this, closed-form expressions are derived for these measures under a cubic distribution and propose new robust quantile moment-based (RQM) measures. These are shown to offer improved empirical performance. The resulting time-varying skewness and kurtosis series are benchmarked against those obtained via maximum likelihood methods. The method provides a practical, interpretable, and robust alternative for modeling higher-order moments directly relevant to risk management and asset allocation.

CO032 Room BCB G09 ML AND HIGH-DIMENSIONAL MACROECONOMIC FORECASTING MODELS
Chair: Anna Simoni
C0405: Measuring and explaining economic uncertainty across US states and metropolitan areas

Presenter: **Eric Ghysels**, University of North Carolina Chapel Hill, United States

Co-authors: Gerald Cohen, Cody Morris, Jieyao Wang

Novel measures are constructed of economic uncertainty at the state and extended metropolitan area (EMA) levels using forecast errors from predictive models of GDP growth. At the state level, uncertainty measures are developed based on sparse-group LASSO MIDAS regressions, incorporating structural adjustments for crisis periods. For EMAs, a stacked sg-LASSO MIDAS fixed effects model is estimated, and a forecast error decomposition is applied to extract quarterly uncertainty shocks from annual GDP data. Results show that EMA uncertainty follows similar patterns to state and macroeconomic uncertainty measures but exhibits distinct fluctuations in normal periods, providing new insights into localized economic uncertainty.

C0573: Factor-augmented sparse MIDAS regressions with an application to nowcasting

Presenter: **Jonas Striaukus**, Copenhagen Business School, Denmark

Co-authors: Jad Beyhum

The purpose is to investigate factor-augmented sparse MIDAS (mixed data sampling) regressions for high-dimensional time series data, which may be observed at different frequencies. The novel approach integrates sparse and dense dimensionality reduction techniques. The convergence rate of the estimator is derived under misspecification due to the MIDAS approximation error, τ -mixing dependence, and polynomial tails. The method's finite sample performance is assessed via Monte Carlo simulations. The methodology is applied to nowcasting U.S. GDP growth, and it is demonstrated that it outperforms both sparse regression and standard factor-augmented regression during the COVID-19 pandemic. These findings indicate that the growth during the pandemic was influenced by both idiosyncratic (sparse) and common (dense) shocks.

C0800: The quantile regression forest reloaded: The distributional random forest

Presenter: **Jeffrey Naef**, University of Geneva, Switzerland

Co-authors: Joan Paredes

Assessing the joint risks of inflation and economic activity, so-called stagflation risks, is a key priority for central banks with medium-term-oriented mandates. A novel approach is introduced based on multivariate evolution of the quantile regression forest, the distributional random

forests (DRF), which allows for the flexible modeling of the joint distribution of inflation and activity, conditional on a rich set of macroeconomic predictors. The method extends traditional random forest techniques by jointly estimating the conditional distribution of multiple outcomes, thereby capturing complex, potentially nonlinear interactions and co-dependencies relevant for stagflation dynamics. Drawing inspiration from the Barro-Gordon framework, the focus is on the probability of scenarios where inflation remains elevated while output contracts, events that are particularly challenging for monetary policy.

C1161: Misspecification, factors, and group sparsity for macro-at-risk models

Presenter: **Matteo Mogliani**, Banque de France, France

Co-authors: Anna Simoni

The aim is to consider a framework where the macroeconomist disposes of a large set of variables for estimating macro-at-risk models. In practice, one can consider either dense or (approximate) sparse models with group structure. These two approaches are analyzed in different scenarios that account for: (i) diverse degrees of sparsity, (ii) signal-to-noise ratios in the probabilistic mechanism generating the covariates, and (iii) alternative correlation structures between groups. A Bayesian quantile regression model is used based on the asymmetric Laplace (AL) distribution and different shrinkage priors. This model can be misspecified in two dimensions: The AL likelihood function and the quantile regression model. It is shown that the first type of misspecification does not impact, in general, the asymptotic results. On the other hand, the second type of misspecification is more serious, and the asymptotic performance depends on the degree of sparsity assumed and the prior adopted. It is shown that there are prior distributions that are more robust than others to misspecification of the macro-at-risk model. Monte Carlo simulations and real data exercises provide further illustration of the results.

CO218 Room Virtual R01 ADVANCED METHODS IN SINGLE-CELL GENOMICS

Chair: Chong Jin

C0994: From single cells to spatial maps: Computational alignment of transcriptomic data

Presenter: **Yunlu Chen**, Northwestern University, United States

Co-authors: Feng Ruan, Ji-Ping Wang

Spatial transcriptomics (ST) measures mRNA transcripts at thousands of locations within tissue slices, revealing spatial variations in gene expression and cell types. However, ST data lacks single-cell resolution, necessitating computational methods to align single-cell RNA-seq (scRNA-seq) data with ST. A novel alignment method is developed with an accompanying Python package, NLSDeconv, based on non-negative least squares for efficient cell-type deconvolution and spatial mapping of ST data. Benchmarking against 18 existing deconvolution methods across various ST datasets demonstrates the approach's competitive statistical performance and superior computational efficiency. Building on this foundation, the method is extended to handle large-scale ST datasets, and the open-source Python package SLSDDeconv is developed. Looking forward, temporal alignment approaches are explored for ST datasets collected across multiple developmental time points.

C1013: Sceptic: Pseudotime analysis for time-series single-cell sequencing and imaging data

Presenter: **Gang Li**, Changping Laboratory, China

Several computational methods have been developed to construct single-cell pseudotime embeddings for extracting the temporal order of transcriptional cell states from time-series scRNA-seq datasets. However, existing methods suffer from low predictive accuracy, and this problem becomes even worse when generalized to other data types such as scATAC-seq or microscopy images. To address this problem, Sceptic is proposed, a support vector machine model for supervised pseudotime analysis. It is demonstrated that Sceptic achieves significantly improved prediction power relative to state-of-the-art methods, and that Sceptic can be applied to a variety of single-cell data types, including single-nucleus image data.

C1056: Single-cell resolution deconvolution of bulk RNA-seq data via functional regression

Presenter: **Chong Jin**, New Jersey Institute of Technology, United States

Co-authors: Xiaotian Mu

The proliferation of single-cell RNA sequencing (scRNA-seq) data has spurred the development of methods to deconvolve bulk RNA-seq tissues. However, many approaches rely on discrete cell type classifications, which can overlook dynamic cell states and create a false dichotomy between changes in cell composition and cell-type-specific expression. Alternatively, strategies that map bulk data to cell subpopulations may suffer from arbitrary parameter choices and "hard" classifications. BUDGIE (Bulk RNA-seq Deconvolution Using Generalized Inference and Estimation), a novel statistical framework, is proposed that deconvolves bulk RNA-seq data at the single-cell level. BUDGIE models bulk expression using a functional linear regression model. The regression coefficients represent cell abundance and are constrained to vary smoothly over a low-dimensional manifold of the scRNA-seq reference atlas. This "soft" regularization avoids rigid cell classification and captures continuous cell state transitions. Within this framework, statistical testing procedures are developed to identify significant cellular subpopulations associated with specific phenotypes. Through benchmarking against existing methods, the performance of BUDGIE is demonstrated, and it is applied to resolve cellular heterogeneity in bulk cancer tissue data.

C1084: Mist: A hierarchical Bayesian framework for detecting differential DNA methylation dynamics in single-cell

Presenter: **Hao Feng**, The University of Texas Health Science Center at Houston, United States

Recent advancements in single-cell DNA methylation (scDNAm) sequencing technologies have enabled the profiling of epigenetic landscapes at unprecedented resolution, offering insights into cellular heterogeneity, differentiation, and evolution. Trajectory inference, which orders cells along pseudotime, allows researchers to track genomics changes across continuous cell states and identify key loci exhibiting differential methylation. However, no methods currently exist to model methylation changes along pseudotime in scDNAm data. The aim is to present a hierarchical Bayesian framework for scDNAm data analysis. The method, named mist (methylation inference for single-cell along trajectory), models stage-specific biological variations, identifies genomic features with significant methylation changes along pseudotime, and performs differential methylation (DM) analysis across phenotypical groups. Simulations demonstrate its superior accuracy in detecting DM genes along pseudotime compared to existing methods. Applied to multi-omics datasets of mouse embryonic development and developing human brain, mist identifies key developmental regulators, whose methylation patterns align with lineage transitions.

CO054 Room Virtual R02 RECENT ADVANCES IN ANALYZING NETWORK DATA (VIRTUAL)

Chair: Tianxi Li

C0981: Inferring diffusion structures of heterogeneous network cascade

Presenter: **Yubai Yuan**, Penn State University, United States

Network cascade refers to diffusion processes in which outcome changes within part of an interconnected population trigger a sequence of changes across the entire network. These cascades are governed by underlying diffusion networks, which are often latent. Inferring such networks is critical for understanding cascade pathways, uncovering Granger causality of interaction mechanisms among individuals, and enabling tasks such as forecasting or maximizing information propagation. A novel double mixture directed graph model is proposed for inferring multi-layer diffusion networks from cascade data. The proposed model represents cascade pathways as a mixture of diffusion networks across different layers, effectively capturing the strong heterogeneity present in real-world cascades. Additionally, the model imposes layer-specific structural constraints, enabling diffusion networks at different layers to capture complementary cascading patterns at the population level. A key advantage of the model is its convex formulation, which allows establishing both statistical and computational guarantees for the resulting diffusion network estimates. Extensive simulation studies are conducted to demonstrate the model's performance in recovering diverse diffusion structures.

C0991: Network goodness-of-fit for the block-model family*Presenter:* **Jingming Wang**, University of Virginia, United States*Co-authors:* Jiashun Jin, Tracy Ke, Jiajun Tang

The block model family is widely used in network modeling and includes four popular models: SBM, DCBM, MMSBM, and DCMM. However, the question of which block model best fits real networks has received limited attention in the literature. The purpose is to introduce a novel approach using cycle count statistics to address the goodness-of-fit for these block models. By leveraging the cycle count statistics and a network fitting scheme, four GoF metrics with parameter-free limiting distributions of $N(0,1)$ are constructed under the assumed models. These GoF-metrics are applied to some frequently-used real networks for comparison. The numerical results suggest that DCMM is particularly promising for modeling undirected networks.

C1293: Modeling longitudinal AD pathology with bi-directional network systems for protein deposition and brain activities*Presenter:* **Zhiling Gu**, Yale University, United States*Co-authors:* Simiao Gao, Tianxi Li, Yize Zhao

The pathology of Alzheimer's disease (AD) remains incompletely understood, particularly with respect to the interplay between tau protein deposition and brain network connectivity. Clarifying this relationship is essential for improving diagnosis, staging, and progression prediction for AD. The aim is to propose a novel framework to characterize the bidirectional interactions between tau accumulation and brain connectivity over time. This framework is motivated by the reciprocal association whereby tau deposition disrupts functional and structural connectivity, while connectivity patterns influence the spread of tau across regions. The approach provides a new avenue for studying the longitudinal co-evolution of protein pathology and brain networks. Through extensive simulations, the consistency of the proposed estimator is demonstrated. Application of data from the Alzheimer's disease neuroimaging initiative further reveals meaningful associations between tau burden and brain connectivity, offering new insights into disease mechanisms.

C1480: Optimal inference on two-sample network effects for directed networks*Presenter:* **Wen Zhou**, New York University, United States*Co-authors:* Yuan Zhang, Mladen Kolar, Yating Liu

Directed networks model a wide range of asymmetric relationships, including social media connections, communication structures, and economic trade flows. These networks are often sparse and exhibit complex edge dependencies such as reciprocity, same-sender, same-receiver, and sender-receiver effects, each capturing key structural patterns like mutual cooperation, broadcasting behavior, recipient diversity, or indirect influence. While prior work has focused on estimating these effects within a single (typically dense) network, understanding how they change across two networks remains critical—particularly for studies of evolving network dynamics or policy interventions. We propose an adaptively unified motif-based subsampling framework for two-sample comparison in directed networks in testing for differences in each of the four edge dependencies. Our approach addresses several core challenges in two-sample network inference, including: (i) differing network sizes and node sets, (ii) varying sparsity levels, (iii) reliance on strong distributional assumptions, and most importantly (iv) the critical issue of unbalanced or indeterminate degeneracy. We derive a Berry-Esseen bound for our test statistics and establish near-optimal performance under various degeneracy scenarios.

CO019 Room BCB 206 CFE SESSION: A TRIBUTE TO H. PESARAN**Chair: Joachim Schnurbus****C0389: High-dimensional causal projection estimators under weak confounding***Presenter:* **Alessio Sancetta**, Royal Holloway, University of London, United Kingdom

Abstract projection estimators are a fundamental tool in applied economics, widely used to approximate causal relationships when the true model is unknown. High-dimensional causal projection estimators that remain robust to unobserved confounding are introduced. A framework is considered in which a large number of observed covariates can offset the bias induced by pervasive, latent confounders. The approach accommodates variation in confounder strength, allowing for weak confounding, and yields estimators that are consistent and asymptotically normal, even under fat-tailed and weakly dependent data. The methodology is demonstrated in three empirical settings: Estimating the dynamic effects of monetary shocks, revisiting GDP convergence across countries, and assessing the impact of U.S. expropriation laws on house prices.

C0701: Monitoring the health of the macroeconomy with a VAR-E dashboard*Presenter:* **Kevin Lee**, University of Nottingham, United Kingdom*Co-authors:* Kalvinder Shields

The time series properties of output, price inflation, and interest rates can be accurately captured using VAR-Es, vector autoregressive models of actual and expected measures of the series, where the latter are provided by surveys. A method is proposed for estimating VAR-Es that accommodates individuals' real-time understanding of the macroeconomy and that delivers forecasts in a way that is useful to decision-makers. A summary is provided of the sort of statistics and figures that might be reported in a dashboard to monitor the health of the macroeconomy, and this is illustrated using the actual and expected data produced by the Bank of England's decision-maker panel.

C0779: Recent advances in the econometric analysis of linear rational expectations models*Presenter:* **Majid Al Sadoon**, Durham University, United Kingdom

The aim is to review recent progress in the specification, solution, identification, and estimation of linear rational expectations models. The aim is to place these developments in historical context, emphasizing Hashem Pesaran's contributions. Particular attention is given to complexity, stationarity, and cointegration. It concludes with a number of conjectures and open problems.

C0502: Univariate properties and the limits to multivariate predictability*Presenter:* **Donald Robertson**, University of Cambridge, United Kingdom*Co-authors:* Stephen Wright

Low-order univariate ARMA processes fit the data nearly as well as multivariate models, and frequently do better at out-of-sample prediction. Multivariate VAR/ABCD models also typically imply much higher-order univariate reduced forms. These observations are reconciled by establishing conditions for multivariate systems to improve upon (fundamental) ARMA. A prior study shows that, for any given series, the maximal-predictive model arises from a univariate representation with hidden (i.e., nonfundamental) univariate shocks. If these representations are treated as the DGP, it is shown that multivariate models can only improve on the fundamental ARMA if a) the nonfundamental univariate shocks are sufficiently correlated and b) the univariate representations of different series have sufficiently different moving average parameters. Empirically, these conditions do not appear to apply to a standard set of macro variables.

CO177 Room BCB 207 TRACKING BUSINESS CYCLES IN HIGH DIMENSIONAL AND/OR VOLATILE SETTINGS**Chair: Andrew Butters****C1147: Monthly GDP growth estimates for the U.S. states***Presenter:* **Aristeidis Raftapostolos**, Kings College London, United Kingdom*Co-authors:* Gary Koop, Stuart McIntyre, James Mitchell

The purpose is to develop a mixed frequency vector autoregressive (MF-VAR) model to produce nowcasts and historical estimates of monthly real state-level GDP for the 50 U.S. states, plus Washington DC, from 1964 through the present day. The MF-VAR model incorporates state and U.S. data at the monthly, quarterly, and annual frequencies. Temporal and cross-sectional constraints are imposed to ensure that the monthly

state-level estimates are consistent with official estimates of quarterly GDP at the U.S. and state-levels. The utility of the historical estimates in better understanding state business cycles and cross-state dependencies is illustrated. It is shown how the model produces accurate nowcasts of state GDP three months ahead of the BEA's quarterly estimates, after conditioning on the latest estimates of U.S. GDP.

C1452: Multi-sector business cycle accounting in a data-rich environment

Presenter: **Andrew Butters**, Indiana University, United States

Co-authors: Scott Brave

Motivated by a multi-sector general equilibrium model with input-output linkages, a structural dynamic factor model is used to decompose U.S. macroeconomic fluctuations into the contributions of shocks to the four "wedges" commonly used in business cycle accounting: (i) an efficiency, (ii) a labor, (iii) an investment, and (iv) a government wedge. The extent to which shocks to these wedges account for the degree of cross-sectional co-movement in a panel of nearly 150 macroeconomic indicators at business cycle frequencies is then evaluated. Evidence is found that the investment and labor wedges are the most likely source of this qualitative feature of business cycles for the U.S., but that specific features of the volatility in the latest business cycle (i.e., pandemic recession) need to be accounted for for these wedges to play an important role recently.

C0464: Nowcasting French GDP with macroeconomic random forests

Presenter: **Julien Andre**, University Paris Dauphine/Banque de France, France

The performance of three types of random forests is compared and evaluated in a recent period. Random forests with and without predictors are reviewed in a setting with mixed frequencies and missing values. The primary focus is on their ability to nowcast French GDP growth during the Covid crisis. Monte Carlo simulations are used to validate estimations, and variable contributions are assessed using variable importance tools.

C0890: Nowcasting GDP using time-varying machine learning methods and mixed-frequency data

Presenter: **Giacomo Caterini**, Italian Parliamentary Budget Office, Italy

Co-authors: Cecilia Frale, Libero Monteforte

Advances in data availability and storage have provided forecasters with an increasingly rich set of economic indicators. Given the delayed and uncertain release of official GDP figures, nowcasting models have emerged to estimate current-quarter activity in real time. The aim is to investigate the time-varying predictive power of large amounts of regressors available with mixed-frequency, exploiting machine learning models for dimensionality reduction and forecasting. Following this approach, it is found that quarterly Italian GDP forecasts based on machine learning algorithms fed by pseudo-real-time monthly information outperform the benchmark, accounting for varying data frequencies and publication lags, addressing the challenges of asynchronous indicator release.

CO416 Room BCB 208 NON-LINEAR IMPULSES AND PROPAGATION MECHANISMS

Chair: Matteo Luciani

C0422: Nonparametric local projections

Presenter: **Elena Pesavento**, Emory University, United States

Co-authors: Silvia Goncalves, Ana Maria Herrera, Lutz Kilian

Nonlinearities play an increasingly important role in applied work when studying the responses of macroeconomic aggregates to policy shocks. Seemingly natural adaptations of the popular local linear projection estimator to nonlinear settings may fail to recover the population responses of interest. The aim is to study the properties of an alternative nonparametric local projection estimator of the conditional and unconditional responses of an outcome variable to an observed identified shock. Alternative ways of implementing this estimator and how to allow for data-dependent tuning parameters are discussed. Results are based on data-generating processes that involve, respectively, nonlinearly transformed regressors, state-dependent coefficients, and nonlinear interactions between shocks and state variables. Monte Carlo simulations show that a local-linear specification of the estimator tends to work well in reasonably large samples and is robust to nonlinearities of unknown form.

C0693: Big vs small monetary policy interventions: A novel approach to nonlinear local projections

Presenter: **Luca Gambetti**, UAB, Spain

Co-authors: Giovanni Ricco, Mario Forni

When the local projection (LP) equations include nonlinear functions of the shock of interest (the nonlinear LP), the impulse response functions estimated using a noisy measure of the shock (instrument or proxy) are distorted. The distortion cannot be eliminated, in general, as in the linear LP, using LP-IV. A solution to this problem is proposed. Under relatively general assumptions, the shock can be estimated as the fitted value of the projection of the instrument onto current and future values of the residuals of a VAR estimated with the set of relevant variables. The estimated shock and its nonlinear functions can then be used as regressors in the LP equations. In the application, it is shown that the effects of US monetary policy shocks on real activity variables increase with the magnitude of the shock, while no size asymmetry is found for prices. The evidence calls for a small tightening to control inflation but a large easing to stimulate the economy.

C0802: Belief distortions and disagreement about inflation

Presenter: **Lorenza Rossi**, University of Lancaster and University of Pavia, United Kingdom

The aim is to investigate the macroeconomic effects of a belief distortion shock and an unexpected increase in the wedge between household and professional forecaster inflation expectations. Using survey and macro data alongside machine-learning techniques, this shock is identified and its effects are examined within and outside the ZLB, while conditioning on the degree of inflation disagreement. The shock increases unemployment during normal times, whereas it reduces it in the ZLB, when the monetary stance is accommodative. Inflation disagreement instead dampens the expansionary effects of the shock. A New Keynesian model with belief distortion shocks replicates these dynamics and reproduces the inflation disagreement empirical patterns.

C0849: Business cycles and financial conditions: A non-parametric vector autoregression (NP-VAR) approach

Presenter: **Francesco Furno**, Amazon.com, United States

Co-authors: Domenico Giannone

The aim is to estimate the joint intertemporal distribution of economic and financial conditions across recessions and expansions in the US. To do so, a non-parametric vector autoregression (NP-VAR) that includes a binary recession variable, business sentiment, and financial conditions is proposed. Results show that recession risk emerges with the simultaneous occurrence of deteriorating economic conditions and tightening financial conditions.

CO155 Room BCB 209 ADVANCES IN MACHINE LEARNING FOR COMMODITY FORECASTING

Chair: Morten Ristad

C0287: Diagnosing cryptocurrency security vulnerability through time series decomposition

Presenter: **Julien Chevallier**, University Paris 8 (LED), France

Cryptocurrency markets have experienced repeated systemic breakdowns over the past decade. Notable examples include the Mt. Gox collapses (2011-2014), the COVID-19-driven "312" flash crash, the "519" crash in 2021 following environmental concerns, the crypto winter of 2017-2018, the collapse of the Terra/Luna ecosystem and the FTX exchange in 2022, as well as the combined impact of Grayscale ETF outflows and tariff conflicts in early 2025. These episodes were triggered by factors such as regulatory shocks, exchange failures, excessive leverage, macroeconomic instability, and automated liquidations, often resulting in billions in losses within hours and trillions in erased market capitalization. Mathematical techniques are proposed for diagnosing such market failures through the decomposition of time series data. The framework integrates nonlinear

signal analysis and regime-shift detection to identify critical transitions before major breakdowns occur. Across security vulnerability episodes, the method highlights consistent patterns in volatility structure, trading flow distortions, and abnormal amplitude shifts preceding failure events. These patterns provide insight into latent instability and offer potential as early diagnostic indicators. A novel lens is contributed for understanding crypto market fragility, offering tools to dissect chaotic behaviors not captured by conventional risk metrics.

C1461: Greening the fleet, raising the price? Shipping decarbonization and inflation dynamics

Presenter: **Soodabeh Ghahramanpour**, Bayes business school, United Kingdom

Co-authors: Ioannis Moutzouris, Malvina Marchese, Ioannis Kyriakou, Mahmoud Fatouh

The shipping industry underpins global trade, carrying over 80% of merchandise worldwide by volume while contributing nearly 3% of global greenhouse gas emissions. As decarbonization policies advance, shipping costs are expected to rise by up to 80% by 2050, potentially influencing inflation in countries highly dependent on maritime transport. This is particularly relevant for the UK, where 85% of traded goods and 40% of imported food rely on shipping. While prior studies confirm a link between shipping costs and inflation levels, their effect on inflation volatility remains underexplored. The aim is to investigate how shipping costs affect UK inflation volatility using a generalized autoregressive conditional heteroscedasticity-in-mean model with exogenous predictors (GARCH-MX). The analysis employs monthly data from 2001 to 2023, combining freight rate indices across major shipping segments (dry bulk, tanker, and container) with key macroeconomic indicators. The results show statistically significant positive effects across all segments, with containerized freight exerting the strongest influence. These findings suggest that shipping costs contribute to inflation uncertainty, particularly as the sector transitions toward net-zero emissions. Policymakers must therefore balance environmental objectives with macroeconomic stability to ensure a sustainable maritime transition without undermining price stability.

C1471: The economics of shipping decarbonization: Carbon, production, and allocative efficiencies

Presenter: **Yao Shi**, City University of London, United Kingdom

Co-authors: Ioannis Moutzouris

The trade-off between environmental and economic performance is investigated in the case of the shipping industry. Existing environmental regulations largely omit the economic efficiency dimension, which, in turn, delays the clean energy transition. A stochastic frontier analysis is applied to assess the relationship between carbon emissions and economic factors as capital, labour, earnings, and transport work, both across all major shipping segments and at an individual-vessel level. The empirical results suggest that technical and operational inefficiencies raise the total cost of owning and operating a vessel by 6%, with market price dynamics and inefficient allocation of economic resources increasing it by 17%. There is scope for the average vessel to reduce its carbon emissions by 31% although carbon efficiency varies significantly depending on the vessel type and period. Higher production efficiency is observed in newer vessels, vessels spending more time at sea, and vessels equipped with more ESTs, and companies with more EST investment. A reduction in speed of up to 15% can improve a vessel's overall efficiency and reduce total cost. Policy interventions need to be carefully designed in order not to negatively impact the overall efficiency of global shipping.

C1213: Monitoring joint tail risks: An application to growth and inflation

Presenter: **Jordi Llorens-Terrazas**, Universidad Carlos III de Madrid, Spain

The aim is to develop the concept of growth and inflation at risk frontier (GIaR). This is a bivariate generalization of the concepts of growth-at-risk (GaR) and inflation-at-risk (IaR). A novel approach is proposed to identify and estimate GIaR and provide uniformly valid upper and lower confidence bands. The procedure is first applied to predict the conditional probability of stagflation. Second, we compute worst-case scenarios for a policy maker who is concerned about the joint tail risk of low growth and high inflation. The effect that a tightening of financial conditions has on the joint tail risks is studied.

CO035 Room BCB 210 RECENT DEVELOPMENTS IN LONG-RUN MODELLING

Chair: Martin Wagner

C1089: IM estimation and (fixed-b) inference for systems of cointegrating multivariate polynomial regressions

Presenter: **Martin Wagner**, University of Klagenfurt and Institute for Advanced Studies, Vienna, Austria

Integrated modified ordinary and generalized least squares estimation are considered for systems of cointegrating multivariate polynomial regressions, i.e., systems of regressions that include deterministic variables, integrated processes, and products of non-negative integer powers of these variables as regressors. The stationary errors are allowed to be correlated across equations, over time, and with the regressors. The necessity to consider integrated modified generalized least squares estimation arises in the case of estimation under restrictions, which in general implies that ordinary and generalized least squares estimators cease to be identical. Hypothesis testing is discussed in detail for the unrestricted and restricted estimators. Furthermore, asymptotically pivotal fixed-b inference is developed, which is shown to be available only in the case of full design for up-to-the-intercept-identical hypotheses tested in all equations in systems with identical regressors in all equations.

C1090: Systems of multi-factor production functions: Modelling the joint behavior of GDP and emissions

Presenter: **Sebastian Veldhuis**, University of Klagenfurt, Austria

Co-authors: Martin Wagner

Integrated modified least squares estimation is applied and (fixed-b) inference to systems of cointegrated (multi-output) multi-factor translog production functions. Translog systems are a special case of systems of cointegrated multivariate polynomial regressions. As usual in the cointegration literature, the regressors are allowed to be endogenous and the errors correlated, both serially and across equations. In the application, one of the outputs, GDP, is a good, whereas the other considered outputs, CO2 and SOX emissions, are bad. The considered input factors include capital, labor, and different specifications of total factor productivity. The developed methodology is applied to aggregate data for about 20 countries using annual data ranging from about 1950 to 2020. Emphasis is put on testing economically relevant hypotheses, e.g., constant returns to scale or the Cobb-Douglas functional form of the GDP equation.

C1329: Revisiting the productivity of public capital

Presenter: **Lennart Empting**, University of Duisburg-Essen, Germany

Co-authors: Helmut Herwartz

Recent policy trends like the Green New Deal draw attention to public investment again, and therewith, to an empirical question of long-lasting presence in academic discussions: Does spending on public capital promote economic growth? As the dynamic effects are still unclear, we revise the empirical literature in this article. Based on the new, extensive data set for 23 OECD countries over the years 1960-2019, vector autoregressive (VAR) methods for heterogeneous panels are employed, and the established approaches of the literature are discussed step by step. It is found that the production function is the single long-run relation in the panel literature, and the varying cointegration ranks of the VAR literature can be summarized by $r = 2$ homogeneous cointegrating vectors in a panel vector error correction model (VECM). Results further confirm that a time-decreasing productivity already found in some VECM studies partially explains the puzzling and often insignificant country-specific impulse response functions (IRF).

CO257: Mispricing proxies in factor models for asset returns

Presenter: **Carlo A Favero**, Bocconi University, Italy

Co-authors: Iliaria Ioni, Gabriele Confalonieri

The aim is to examine the influence of mispricing proxies on stock return dynamics within the framework of Fama-French five-factor models. Specifically, the role of mispricing proxies derived from cointegration is assessed between asset prices and factor prices, as well as sentiment

indicators extracted from quarterly earnings conference calls. Using quarterly data from 1980 to 2023 for the cross-section of DJIA-listed firms, the empirical analysis shows that deviations from long-run trends, driven by factor prices, have predictive power for stock returns after controlling for the five Fama-French factors. Stock-specific sentiment further enhances predictability. The additional predictability generated by mispricing proxies is fully explained by a non-linear model where sentiment determines the speed of adjustment toward the long-run trend identified by cointegration analysis when stock prices are above it.

CO331 Room BCB 211 MODELLING TIME-VARYING RELATIONSHIPS IN ECONOMETRICS
Chair: Liudas Giraitis
C0570: Strong spatial dependence

Presenter: **Jungyoon Lee**, Royal Holloway, University of London, United Kingdom

Co-authors: Francesca Rossi, Offer Lieberman

The purpose is to study the asymptotic behavior of the two-stage least squares estimation in spatial autoregression when the spatial dependence is strong. Motivated by periods of exuberance in the housing market leading to estimates of spatial parameters in the explosive range, the asymptotic properties of estimates in such models are explored, linking the results to those available in the time series literature.

C0628: Augmented dynamic model

Presenter: **Yufei Li**, Kings College London, United Kingdom

Co-authors: Liudas Giraitis, George Kapetanios

The recent work on regression modeling that permits general heterogeneity is extended to allow for lagged dependent variables. The purpose is to explore to what extent the generality of the setting, the simplicity of assumptions, and the ease of computation of standard errors can be preserved. Theoretical properties of regression estimation and inference is accompanied by Monte Carlo experiments and an empirical application.

C0710: Double-pooling for dynamic tail estimation

Presenter: **Shiqi Ye**, AMSS Center for Forecasting Science, Chinese Academy of Sciences, China

The aim is to propose a novel semiparametric framework for modeling the dynamic tail behavior of panel time series. The methodology introduces a "double-pooling" mechanism that effectively exploits both cross-sectional dependence and cross-quantile structures. This approach significantly enhances the accuracy and efficiency of estimating time-varying value-at-risk (VaR) and expected shortfall (ES). The asymptotic theory for the estimators is established, and formal hypothesis tests are developed to rigorously evaluate the proposed tail dynamics. An empirical application to European economies demonstrates the practical utility of the method in providing robust insights into dynamic growth-at-risk and expected shortfall, offering valuable implications for financial stability and risk management.

C0162: Beyond the mean: Limit theory and testing in infinite-mean acd duration models

Presenter: **Anders Rahbek**, University of Copenhagen, Denmark

Co-authors: Giuseppe Cavaliere, Thomas Mikosch, Frederik Vilandt

Integrated autoregressive conditional duration (ACD) models serve as natural counterparts to the integrated GARCH models used for financial returns. Despite their resemblance, asymptotic theory for ACD is challenging and also not complete, in particular for integrated ACD. Central challenges arise from the facts that (i) integrated ACD processes imply durations with infinite expectation, and (ii) even in the non-integrated case, conventional asymptotic approaches break down due to the randomness in the number of durations within a fixed observation period. Addressing these challenges, we provide unified asymptotic theory for the (quasi-) maximum likelihood estimator for ACD models; a unified theory which includes integrated ACD models. We also provide a novel framework for hypothesis testing in duration models, enabling inference on a key empirical question: whether durations possess a finite or infinite expectation. We apply our results to high-frequency cryptocurrency ETF trading data. Motivated by parameter estimates near the integrated ACD boundary, we assess whether durations between trades in these markets have finite expectation, an assumption often made implicitly in the literature on point process models. Our empirical findings indicate infinite-mean durations for all the five cryptocurrencies examined, with the integrated ACD hypothesis rejected against alternatives with tail index less than one for four out of the five cryptocurrencies considered.

CO024 Room BCB 212 NON-GAUSSIAN TIME SERIES
Chair: Joann Jasiak
C0538: Intraday functional PCA forecasting of cryptocurrency returns

Presenter: **Cheng Zhong**, York University, Canada

The functional PCA (FPCA) forecasting method is studied in application to functions of intraday returns on Bitcoin. It is shown that improved interval forecasts of future return functions are obtained when the conditional heteroscedasticity of return functions is taken into account. The Karhunen-Loeve (KL) dynamic factor model is introduced to bridge the functional and discrete time dynamic models. It offers a convenient framework for conditionally heteroscedastic-functional time series analysis. For intraday forecasting, a new algorithm is introduced based on the FPCA applied by rolling, which can be used for any data observed continuously 24/7. The proposed FPCA forecasting methods are applied to return functions computed from data sampled hourly and at 15-minute intervals. Next, the functional forecasts evaluated at discrete points in time are compared with the forecasts based on other methods, including machine learning and a traditional ARMA model. The proposed FPCA-based methods perform well in terms of forecast accuracy and outperform competitors in terms of the directional (sign) of return forecasts at fixed points in time.

C0543: Generalized covariance estimator under misspecification and constraints

Presenter: **Aryan Manafi Neyazi**, York University, Canada

The asymptotic properties of the generalized covariance (GCov) estimator are investigated under misspecification. We show that GCov is consistent and has an asymptotically normal distribution under misspecification. Then, GCov-based Wald-type and score-type tests are constructed, all of which follow a chi-square distribution. Furthermore, the indirect GCov (IGCov) and the constrained GCov (CGCov) estimators are proposed. The IGCov estimator is useful for estimating models indirectly, based on simulations, such as non-invertible moving average models. Consequently, an IGCov specification test is developed. The CGCov estimator extends the use of the GCov estimator to a broader range of models with constraints on their parameters. The asymptotic distribution of the CGCov estimator is investigated when the true parameters are far from the boundary and on the boundary of the parameter space. The finite sample performance of proposed estimators and tests is validated in the context of noncausal-noninvertible and DAR models. Finally, two empirical applications are provided by applying the noncausal model to the final energy demand commodity index and also the DAR model to the US 3-month treasury bill.

C0652: Modelling high dimensional realized correlations

Presenter: **Roxana Halbleib**, University of Freiburg, Germany

Co-authors: Ilya Archakov, Jasper Rennspies

A novel approach is proposed to model and forecast high-dimensional realized correlation matrices by using the matrix logarithm transformation introduced by a prior study. To reduce dimensionality, the focus is on extracting factors from the transformed correlation series, i.e., gamma series, according to a particular sector or industry. The empirical results show that the new approach is easily implementable for large-dimensional matrices and outperforms standard methods.

C0990: Nonlinear impulse response functions and local projections*Presenter:* **Christian Gourieroux**, University of Toronto and CREST, Canada*Co-authors:* Quinlan Lee

The goal is to extend the non-parametric estimation of impulse response functions (IRF) by means of local projections in the nonlinear dynamic framework. The existence of a nonlinear autoregressive representation for Markov processes is discussed, and it is explained how their IRFs are directly linked to the nonlinear local projection (NLP), as in the case for the linear setting. A fully non-parametric LP estimator is presented in the one-dimensional nonlinear framework, and its asymptotic properties are compared to those of IRFs implied by the autoregressive model. Extensions are also considered to the multivariate framework through the lens of semiparametric models.

CO101 Room BCB 213 RECENT ADVANCES IN QUANTILE REGRESSION**Chair: Carlos Lamarche****C0668: Model averaging in semiparametric estimation of quantile treatment effects***Presenter:* **Antonio Galvao**, Michigan State University, United States*Co-authors:* Sergio Firpo, Ulrich Hounyo, Li Lu

Model-averaging methods are proposed to estimate quantile treatment effects (QTE and QTT) under treatment selection based on observables. To address propensity score misspecification, two estimators are developed: One averaging QTE/QTT across models and another that averages propensity scores before estimation. Unconfoundedness is required for at least one covariate set or their union. A data-driven covariate selection criterion is introduced, and asymptotic properties are derived for inference. A novel "unconfoundedness signature plot" is introduced and helps to assess covariate relevance. Simulations show strong finite-sample performance, and the approach is illustrated by estimating the effect of inherited control on firm performance.

C0709: Parametric bootstrap inference for quantile regression and binary quantile regression*Presenter:* **Blaise Melly**, University of Bern, Switzerland*Co-authors:* Martina Pons

The purpose is to develop (semi)parametric bootstrap inference methods for quantile and binary quantile regression models. The complete quantile regression process characterizes the entire conditional distribution of the outcome given the covariates, allowing outcomes to be resampled directly from the estimated distribution. Moreover, imposing a null hypothesis, such as a location shift, during the simulation of bootstrap samples is straightforward. This approach yields improved size and power compared to nonparametric bootstrap methods that rely on ex post recentering. The benefits of this parametric bootstrap are even greater in the context of binary quantile regression (generalized maximum score), where inference is particularly challenging due to the estimator's nonstandard asymptotic distribution and the failure of standard bootstrap techniques. By simulating outcomes from the estimated binary quantile regression process, the method provides a practical and effective approach to valid inference without requiring the selection of smoothing parameters.

C0325: Quantile regression inference processes and choices in sparsity estimation*Presenter:* **Thomas Parker**, University of Waterloo, Canada

The aim is to investigate uniform inference for conditional quantile functions, in the spirit of a prior study. The focus is on estimation of the derivative of the conditional quantile process, or the sparsity function, beyond the choices proposed by prior studies. Recent technical advances have proposed methods to establish the uniform consistency of density estimates that can be adapted to this setting, allowing researchers to use a wide variety of sparsity estimators for inference. A small simulation experiment compares the finite sample performance of a few such estimators with the well-established sparsity estimates.

C0566: Partitioned wild bootstrap for panel data quantile regression*Presenter:* **Carlos Lamarche**, University of Kentucky, United States*Co-authors:* Antonio Galvao, Thomas Parker

Practical inference procedures for quantile regression models of panel data have been a pervasive concern in empirical work, and can be especially challenging when the panel is observed over many time periods and temporal dependence needs to be taken into account. A new bootstrap method is proposed that applies random weighting to a partition of the data (partition-invariant weights are used in the bootstrap data-generating process) to conduct statistical inference for conditional quantiles in panel data that have significant time-series dependence. It is demonstrated that the procedure is asymptotically valid for approximating the distribution of the fixed effects quantile regression estimator. The bootstrap procedure offers a viable alternative to existing resampling methods. Simulation studies show numerical evidence that the novel approach has accurate small-sample behavior, and an empirical application illustrates its use.

CO038 Room BCB M201 STATISTICAL INFERENCE AND PREDICTION OF EXTREME VALUES AND BEYOND**Chair: Simone Padoan****C0366: Statistical prediction of peaks over a threshold***Presenter:* **Stefano Rizzelli**, University of Padova, Italy*Co-authors:* Simone Padoan

In many applied fields, the prediction of more severe events than those already recorded is crucial for safeguarding against potential future calamities. What-if analyses, which evaluate hypothetical scenarios, play a key role in assessing the potential impacts of extreme events and in guiding the development of effective safety policies. These problems can be tackled using extreme value theory. The well-established peaks-over-threshold method is employed, and a comprehensive toolkit is described to address forecasting needs. An out-of-sample variable is examined, and the focus is on its conditional probability of exceeding a high threshold. Conditions are given under which the generalized Pareto approximation of the corresponding predictive density is accurate, and a Bayesian approach is described for its estimation, enabling the derivation of predictive intervals. By leveraging threshold stability, it is illustrated how predictions can be reliably extended deep into the tail of the unknown data distribution. Asymptotic accuracy of the proposed estimator and predictive intervals is established, as well as that of estimators of notable risk measures based on point forecasts. Finally, the prediction framework is extended to the case of independent data with covariates within a proportional tail model, and to the case of linear time series.

C0586: Neural statistical modelling of cascading extremes*Presenter:* **Miguel de Carvalho**, University of Edinburgh and Universidade de Aveiro, Portugal*Co-authors:* Clemente Ferrer, Ronny Vallejos

The purpose is to address the growing concern surrounding cascading extreme events, such as a major undersea earthquake triggering a tsunami, by presenting a novel method for risk assessment focused on these domino effects. The proposed approach develops an extreme value theory framework within a Kolmogorov-Arnold network (KAN) to estimate the probability of one extreme event triggering another, conditionally on a feature vector. An extra layer is added to the KANs architecture to enforce the definition of the parameter of interest within the unit interval, and the resulting neural model is referred to as KANE (KAN with natural enforcement). The method is supported by extensive numerical studies and further demonstrated through real-world applications. Finally, connections between the proposed methods and ongoing work on generative modeling of extreme events are mentioned.

C0677: Nonparametric heavy-tailed distribution estimation via random probability measures*Presenter:* **Carlotta Pacifici**, Bocconi University, Italy*Co-authors:* Simone Padoan, Stefano Rizzelli

Estimating the distribution that has generated the observed data is particularly challenging in the case of heavy tails. This task is addressed by focusing on three key objectives: Avoiding strong modeling assumptions, accurately capturing the upper-tail behavior, and allowing for uncertainty quantification. A Bayesian nonparametric framework is adopted, using the Pitman-Yor process (PYP) for estimating the unknown data-generating distribution. Unlike the Dirichlet process, the PYP preserves heavy tails when centered on a heavy-tailed base measure, the same holds for the posterior process. It is sampled from the posterior PYP that combines information from the observed data and the prior base measure, resulting in multiple trajectories that aim to mimic the sample distribution. Simulation studies show coherence between the sampled trajectories and the true one when the base measure matches the true model. Since the data-generating process is unknown in practice, as base measure, a piecewise density is specified with a generalized Pareto upper tail. The resulting trajectories capture the true tail index and yield good coverage at both low and high portions of the distribution.

C0942: Extreme-value modelling of migratory bird arrival dates: Insights from citizen science data*Presenter:* **Jonathan Koh**, ETH Zurich, Switzerland

Citizen science mobilizes many observers and gathers huge datasets, but often without strict sampling protocols, resulting in observation biases due to heterogeneous sampling effort, which can lead to biased predictions. A spatio-temporal Bayesian hierarchical model is developed for bias-corrected estimation of arrival dates of the first migratory bird individuals at their breeding sites. Higher sampling effort could be correlated with earlier observed dates. Data fusion of two citizen-science datasets is implemented with fundamentally different protocols (Breeding Bird Survey, eBird), and posterior distributions of the latent process are obtained, which contain four spatial components endowed with Gaussian process priors: Species niche, sampling effort, position, and scale parameters of annual first arrival date. The data layer consists of four response variables: Counts of observed eBird locations (Poisson), presence or absence at observed eBird locations (Binomial), BBS occurrence counts (Poisson), and first arrival dates (generalized extreme-value). The aim is to devise a Markov chain Monte Carlo scheme and check by simulation that the latent process components are identifiable. The model is applied to several migratory bird species in the northeastern US for 2001-2021, and it is found that the sampling effort significantly modulates the observed first arrival dates. This relationship is exploited to effectively bias-correct predictions of the true first arrivals.

CO192 Room BCB M202 THE FINANCIAL RISKS OF CLIMATE CHANGE AND BIODIVERSITY LOSS Chair: Juan-Angel Jimenez-Martin
C0298: Evaluating clean energy's impact on forecasting probability distribution in the energy market*Presenter:* **Laura Garcia-Jorcano**, Universidad de Castilla-La Mancha, Spain*Co-authors:* Lidia Sanchis-Marco

Loss distribution prediction is investigated in the energy sector of the S&P 500 using two distinct modeling approaches: The standard SAV-CAViaR model and an extended version that incorporates the Wilder-Hill clean energy index (ECO) as a proxy for the clean energy transition. Through an analysis of the distributions of daily returns for the standard and extended models, the aim is to evaluate these models' predictive power across various quantile levels ranging from 0.01 to 0.99. A comprehensive analysis of the obtained distributions is provided, including statistical moments, and assesses the out-of-sample predictive performance of the clean energy variable. New insights are offered into the impact of clean energy transitions on financial risk, underscoring the importance of integrating clean energy variables into financial risk models. It demonstrates that these variables significantly impact predictive accuracy and inform sustainable investment strategies.

C0302: Extreme climate and natural disaster risk in financial markets: A CoES approach*Presenter:* **Lidia Sanchis-Marco**, University of Castilla-La Mancha, Spain*Co-authors:* Laura Garcia-Jorcano

The financial impact of extreme climate and natural disaster (EC&ND) events is analyzed for S&P 500 sectors, focusing on their role in increasing market and systemic risk. Using the conditional autoregressive expected shortfall (CARES) model, sector-specific losses are estimated and forecasted by incorporating EC&ND variables. An enhanced conditional expected shortfall model is also introduced to assess systemic risk through the dependence between financial sector losses and overall system losses. Compared to traditional models such as value-at-risk (VaR) and conditional VaR, this approach provides a more detailed understanding of tail risk under climate stress. Results show that EC&ND variables significantly increase estimated and forecasted losses, especially in energy, real estate, and insurance, the sectors most sensitive to climate extremes. Two new metrics, the disaster market ratio and the systemic risk ratio, are also proposed to quantify these effects. Overall, findings emphasize the urgent need to integrate EC&ND risk into financial regulation and policy frameworks to strengthen systemic resilience. By linking climate risk to financial stability, practical insights are provided for regulators, investors, and academics.

C0828: Do the elements of biodiversity affect financial risk? An analysis for US firms.*Presenter:* **Almudena Maria Garcia Sanz**, Complutense University of Madrid, Spain*Co-authors:* M Dolores Robles, Juan-Angel Jimenez-Martin

The aim is to examine the role of biodiversity exposure and the financial risk of firms between 2010 and 2022. A set of US companies included in the Refinitiv DataStream indexes is analyzed. A panel regression analysis is applied to explain a comprehensive set of firms' risk dimensions as a function of a set of recent biodiversity metrics developed in a recent study, as well as other metrics obtained from Thomson Reuters. Results indicate that biodiversity impacts financial risks in different ways, even after controlling for firm characteristics such as size and profitability. These findings offer useful insights for all stakeholders, especially from the policy-makers' perspective to know more clearly how to impulse more diverse and sustainable economies, and from the academic perspective too, as they improve the understanding of the impact of companies' ESG and biodiversity preservation engagement on their performance.

C1403: Semantic weighting of climate solutions and their impact on sustainable profitability of firms*Presenter:* **Yanchen Liu**, Universidad Complutense de Madrid, Spain

Using a semantically weighted measure of climate solutions constructed from the business description sections of the 10-K filings of US public firms, its relationship with sustainable profitability is examined. Unlike frequency-based indicators that count the proportion of climate solution sentences but fail to distinguish between substantive and superficial mentions, our approach applies a deep contextualized term weighting model to assign relevance scores, producing an indicator that better captures the quality of climate solution activities. The measure shows stronger correlation with external environmental innovation scores and superior explanatory power in panel regressions relative to the standard proportional indicator. Analyzing profitability outcomes, it is found that high scores on this measure are associated with lower gross margins, higher returns on sales, and a dual effect on return on assets, reflecting both efficiency premiums and overinvestment penalties. Overall, results suggest that semantic weighting enhances the measurement of climate solutions and reveals the nuanced financial implications of climate-related strategies of firms.

CO090 Room BCB 308 STATISTICAL LEARNING WITH MIXTURES
Chair: Alejandro Murua
C0950: Model-based color quantization with Gaussian copula mixtures*Presenter:* **Ranjan Maitra**, Iowa State University, United States*Co-authors:* Samipan Majumder

Color quantization is used in computer graphics to reduce the number of colors in an image without appreciably losing its visual quality. The importance of this process comes from a need to display images on devices that are not completely capable of displaying all the colors in a digital image. Each pixel, or picture element, in an image is represented in terms of its primary components, namely Red, Green, and Blue. Therefore, each pixel has a certain amount of Red, a certain amount of Green, and a certain amount of Blue. This way of representing color is known as the RGB format. Model-based color quantization is performed in a way that more accurately models the RGB channels of a digital image. The model uses a Gaussian copula with Kumaraswamy marginals, while also practically incorporating spatial context through penalization. The methodology is demonstrated on several digital images.

C0954: Learning sparse mixture-of-experts generalized linear models in ultrahigh dimensions

Presenter: **Pengqi Liu**, McGill University, Canada

Co-authors: Abbas Khalili

Mixture-of-experts generalized linear models (MoE-GLM) are used for analyzing data that arise from populations with unobserved heterogeneity. In recent applications of MoE-GLM, data are often collected on a large number of features. However, fitting an MoE-GLM to such high-dimensional data is numerically challenging. To cope with the high-dimensionality in estimation, it is often assumed that the model is sparse and only a handful of features are relevant to the analysis. Most of the existing development on sparse estimation is in the context of homogeneous regression or supervised learning problems. The focus is on some of the challenges and also recent computational and theoretical developments for sparse estimation in MoE-GLM when the number of features can be in exponential order of the sample size (ultrahigh-dimensional setting). The asymptotic properties of the proposed methodology and its performance in finite-sample settings are also presented and discussed via simulations and a real data analysis.

C0939: Generalized multilinear models for tensor-on-tensor regression

Presenter: **Carlos Llosa**, Sandia National Laboratories, United States

Co-authors: Daniel M Dunlavy, Richard B Lehoucq, Jeremy Myers, Tian Ma

The generalized multilinear model (GMLM) is introduced, a novel modeling framework that extends the generalized linear model (GLM) and tensor-on-tensor regression (ToTR) for regression problems involving tensor, or multidimensional array, data. As with GLMs, GMLMs allow a linear model to relate expected response variables following arbitrary distributions to covariates via a general link function, providing flexibility in solving problems beyond typical identity link/Gaussian-response regression. As in ToTR, GMLMs allow for tensor covariates and responses, providing models that can leverage the multilinear structure inherent in many data that is often discarded when the data is vectorized and modeled entry-wise using scalar-response GLMs. Vectorizing the data often leads to an ill-posed problem unless provided a large sample size that increases with the product of the sizes of the covariate and response tensors. Instead, a low-rank tensor structure is imposed on the GMLM parameter tensor, thus requiring fewer samples and leading to a well-posed inference problem. The extensions of GLMs and ToTR that lead to GMLMs are discussed, an algorithmic framework is introduced for solving the GMLM parameter inference problem when the low-rank structure imposed on the parameter tensor is the Canonical Polyadic (CP) model, and multiple uses of GMLMs are illustrated on simulated and real-world application data.

C1271: A non-parametric integrative Bayesian approach for variable selection and prediction

Presenter: **Thierry Chekouo**, University of Minnesota, United States

Linear models may not adequately capture complex, nonlinear associations between outcomes and features. Moreover, with advancements in technology, data from multiple platforms is now collected on the same individuals, necessitating methods that effectively integrate multi-platform data. To address these limitations, a nonparametric Bayesian variable selection approach that employs Gaussian process priors is proposed to flexibly model the response surface within a model-based data integration framework. The novelties of the method lie in integrating multi-view data using multiple kernels in a Bayesian framework, allowing for simultaneous variable selection for each data type and accurate prediction. The proposed model measures the importance of each data view in predicting clinical outcomes while performing view-specific variable selection. A key feature of the method is that it can also be utilized in accelerated failure time (AFT) models when dealing with censored time-to-event outcomes. Several simulation studies are presented where we demonstrate the capability of the approach to detect significant variables across various data platforms and to predict outcomes effectively.

CO245 Room BCB 309 STARSTRUCK STATISTICS

Chair: Radu Craiu

C0783: Normalizing flows for posterior estimation under intractable likelihoods, with applications in astrostatistics

Presenter: **Roxana Darvishi**, Simon Fraser University, Canada

Co-authors: David Stenning, Owen Ward

Normalizing flows are a class of models that allow for flexible density estimation and efficient sampling of complex distributions. These models use neural networks to learn a transformation that maps a simple random variable to the target of interest. In the context of Bayesian hierarchical models, normalizing flows are particularly useful when the likelihood function is unavailable or computationally expensive to evaluate, limiting the applicability of conventional inference techniques. A two-stage algorithm using normalizing flows is proposed to enable density evaluation and sample generation when the target distribution can be factorized into components that are either available in closed form or accessible through sampling. The method is applied to a challenging astrophysics analysis for which case-by-case posterior samples for several physical parameters of thick-disk white dwarf stars are available, but the explicit mathematical forms of the underlying densities are unknown. The goal is to approximate the full joint posterior and enable both sampling and density evaluation simultaneously. The effectiveness of the algorithm is evaluated by comparing it to existing methods, highlighting its potential advantages for dealing with intractable likelihoods. Broader applications are also discussed, including joint density estimation and reducing the computational cost of sample generation.

C0931: Discovering extremely faint galaxies using MATHPOP

Presenter: **Dayi Li**, University of Toronto, Canada

The aim is to present MATHPOP: A new, mark-dependent thinned point process that infers the number of old star clusters (globular clusters, or GCs) around ultra-diffuse galaxies (UDGs) and low-surface brightness galaxies. Many UDGs have unusually high GC numbers relative to their surface brightness, but standard GC-count methods struggle with photometric uncertainties, membership ambiguities, and assumptions about the GC luminosity function (GCLF). MATHPOP jointly models the spatial and brightness distributions of GCs while minimizing assumptions. MATHPOP is validated against traditional methods using 40 known UDGs and low-surface brightness galaxies in the Perseus galaxy cluster, and has also discovered two intriguing galaxies that appear to have much brighter than average GC populations.

C0978: Astro-statistical learning and uncertainty quantification

Presenter: **Joshua Speagle**, University of Toronto, Canada

Astrophysics is entering a regime where the theoretical models of many astrophysical phenomena are noticeably less expressive than the rich and massive amounts of data we are collecting. As a result, there is increasing interest in leveraging statistical learning methods to simultaneously uncover underlying structure while performing latent parameter inference. The aim is to discuss some of the challenges in accomplishing these goals within astrophysical datasets, which span data characteristics (heterogeneous measurements, biased and censored sampling, etc.), uncertainty

quantification (calibration, covariate shift, etc.), and model selection (ill-posed model classes, etc.), among others. Examples are motivated by particular astrophysical inference problems, including ongoing efforts towards practical solutions.

C1026: Quantifying the clustering probability in noisy nonhomogeneous spatial data to identify repeating fast radio bursts

Presenter: **Amanda M Cook**, McGill University, Canada

The aim is to introduce an approach to analyze nonhomogeneous Poisson processes (NHPP) observed with noise, focusing on previously unstudied second-order characteristics of the noisy process. Utilizing a hierarchical Bayesian model with noisy data, hyperparameters governing a physically motivated NHPP intensity are estimated. Simulation studies are performed to demonstrate the reliability of this methodology in accurately estimating hyperparameters. Leveraging the posterior distribution, the probability of detecting a certain number of events within a given radius, the k -contact distance are then inferred. The methodology is demonstrated with an application to observations of fast radio bursts (FRBs) detected by the Canadian Hydrogen Intensity Mapping Experiment's FRB Project (CHIME/FRB). This approach allows identifying repeating FRB sources by bounding or directly simulating the probability of observing k physically independent sources within some radius, or the probability of chance coincidence (P_{cc}). The new methodology improves the repeater detection P_{cc} , in 86% of cases when applied to the largest sample of previously classified observations, with a median improvement factor (existing metric over P_{cc} , from the methodology) of approximately 3000.

CO083 Room BCB 310 ADVANCES IN LONGITUDINAL AND REPEATED MEASURES DATA

Chair: Sanjoy Sinha

C0326: Advancements in efficient estimation for mixed effects models with censored data

Presenter: **Shakhawat Hossain**, University of Winnipeg, Canada

Longitudinal and repeated measures data are frequently analyzed using mixed models. However, the presence of censored responses, a common occurrence in biomedical research and clinical trials due to detection limits, introduces complexity. These limits can result in left- or right-censored measurements. While adapted linear mixed effects models are often employed to handle such data, a likelihood-based approach is proposed for fitting a linear mixed effects model with normally distributed errors. An expectation-maximization algorithm is utilized for unrestricted maximum likelihood estimation. Furthermore, scenarios where model parameters are subject to uncertain linear constraints are explored, leading to the development of a restricted estimator. To improve the estimation of fixed effects, two refined estimator sets are introduced: A pretest estimator, and shrinkage and positive shrinkage estimators. The performance of these proposed estimators is evaluated against the unrestricted maximum likelihood estimator via extensive simulations and an application to longitudinal data from the AIDS Clinical Trials Group protocol A5055 study.

C0515: Sequential designs for possibly misspecified mixed ANCOVA models with longitudinal data

Presenter: **Xiaojuan Xu**, Brock University, Canada

Co-authors: Sanjoy Sinha

The construction of robust sequential designs are studied for linear mixed models applied to longitudinal or clustered data, where both treatment and covariate effects are present. The objective is to determine optimal allocations of sample units to treatments, as well as appropriate covariate levels, while accounting for possible misspecification in the assumed model structure. The design approach balances two goals: Minimizing the determinant of the estimated variance-covariance matrix to improve precision, and reducing potential estimation/prediction bias arising from model misspecification. The empirical performance of the proposed designs is also examined through simulation studies.

C0999: A goodness-of-fit test for zero-inflated repeated measures: Tree abundance applications

Presenter: **Juxin Liu**, University of Saskatchewan, Canada

Co-authors: Yanyuan Ma, Jill Johnstone

Field studies in ecology often make use of data collected in a hierarchical fashion and may combine studies that vary in sampling design. For example, studies of tree recruitment after disturbance may use counts of individual seedlings from plots that vary in spatial arrangement and sampling density. To account for the multi-level design and the fact that more than a few plots usually yield no individuals, a mixed effects zero-inflated Poisson model is often adopted. Although it is a convenient modeling strategy, various aspects of the model could be misspecified. A comprehensive test procedure, based on the cumulative sum of the residuals, is proposed. The test is proven to be consistent, and its convergence properties are established as well. The application of the proposed test is illustrated by a real data example and simulation studies.

C0869: Constrained inference for longitudinal data with nonignorable missing responses

Presenter: **Sanjoy Sinha**, Carleton University, Canada

Missing data are common in longitudinal studies. When data are nonignorably missing, it is necessary to incorporate the missing data mechanism into the likelihood function to ensure valid inference. The unrestricted maximum likelihood method for incomplete longitudinal data has been extensively studied in the literature. However, parameter orderings or constraints may naturally arise in real-life scenarios, where the efficiency of an estimator can be improved by incorporating these parameter constraints into estimation and hypothesis testing. The purpose is to discuss some novel methods for analyzing longitudinal data with nonignorable missing responses under linear inequality constraints. The proposed method is developed within the framework of the generalized linear mixed model. The empirical properties of the estimators are investigated through Monte Carlo simulations. An application is presented using real data from a health survey.

CO163 Room BCB 312 RECENT ADVANCES IN THE DESIGN OF EXPERIMENTS

Chair: Vasiliki Koutra

C0437: Design of experiments on sampled networks

Presenter: **Luke Ayres**, King's College London, United Kingdom

Co-authors: Vasiliki Koutra

Experiments on social networks are increasingly important, for example, marketing experiments where the effectiveness of different advertisements given to different users needs to be assessed. Development of methods for optimizing the design of experiments on networks is an active area of research. A significant complication is that in networked experiments, the response of a given unit depends not only on the direct treatment applied to that unit, but also on the indirect effect of treatments applied to connected units. Previous research has focused on the problem of optimal design (treatment allocation) to assess direct treatment effects, indirect network effects, or a combination of both. The focus is on how different network properties, such as edge density, impact different optimality criteria. Studying such properties is particularly important for experiments on large networks where it is likely that not all available units will be used due to cost and/or computation time. Hence, a sub-network will need to be chosen for experimentation, with different choices giving different network properties. The use of different network sampling algorithms is assessed to evaluate the effectiveness of the resulting optimal designs and to demonstrate the important role of network structure in determining design efficiency.

C0485: A column generation approach to exact experimental design

Presenter: **Selin Ahipasaoglu**, University of Southampton, United Kingdom

Co-authors: Stefano Cipolla, Jacek Gondzio

The aim is to address the exact D-optimal experimental design problem by proposing an efficient algorithm that rapidly identifies the support of its continuous relaxation. The method leverages a column generation framework to solve such a continuous relaxation, where each restricted master problem is tackled using a primal-dual interior-point-based semidefinite programming solver. This enables fast and reliable detection of the design's support. The identified support is subsequently used to construct a feasible exact design that is provably close to optimal. It is shown

that, for large-scale instances in which the number of regression points exceeds by far the number of experiments, the approach achieves superior performance compared to existing branch-and-bound-based algorithms in both computational efficiency and solution quality.

C0817: **Experimental design and data-driven optimization in complex systems**

Presenter: **Matteo Borrotti**, University of Milan-Bicocca, Italy

Design of experiments (DoE) has long been a cornerstone in planning efficient data collection and enhancing model-based decision-making. However, modern complex systems, characterized by high dimensionality, costly evaluations, and multiple competing objectives, require extending the classical DoE toolbox with new computational and data-driven approaches. The foundations of DoE are revisited, and it is explored how emerging tools, such as Pareto-based selection, advanced optimization criteria, and algorithmic search methods, can be integrated into experimental design strategies. Through illustrative examples and recent developments, it is shown how these methods can support the design of experiments in challenging scenarios, from industrial processes to simulation-based optimization. Finally, open challenges and promising directions where data-driven optimization can enhance the role of DoE in the analysis of complex systems are discussed.

C0866: **Sequential Bayesian design via Laplace policies**

Presenter: **Tim Waite**, University of Manchester, United Kingdom

Co-authors: Emma Rowlinson

Policy-based methods for sequential Bayesian experimental design aim to learn a mapping from the current knowledge state to optimal future experiments, maximizing a utility such as expected information gain. A key consideration is how to represent the knowledge state. The proposal is to use the Laplace approximation to the posterior as this representation, enabling efficient training of neural network policies. The technical foundations of the Laplace policy framework are introduced, and its performance is illustrated across a range of design problems. To enable gradient-based optimization, differentiable computation of the posterior mode is ensured via custom gradient methods based on the implicit function theorem. The framework is further extended to binary response models using concrete relaxation, which enables approximate simulation of discrete random variables while preserving the differentiability of the widely used reparameterization technique for simulation of continuous random variables. New results are presented, including comparisons with state-of-the-art methods such as deep adaptive design (DAD) and policy gradient sequential optimal design (PG-SOED). These highlight the effectiveness of Laplace policies, particularly in settings where the posterior is approximately Gaussian.

CO136 Room BCB 313 MODERN DIMENSION REDUCTION TECHNIQUES AND THEORIES

Chair: Kyongwon Kim

C1188: **Learning causal graphs via nonlinear sufficient dimension reduction**

Presenter: **Kyongwon Kim**, Yonsei University, Korea, South

Co-authors: Eftychia Solea, Bing Li

The purpose is to introduce a new nonparametric methodology for estimating a directed acyclic graph (DAG) from observational data. The method is nonparametric in nature: It does not impose any specific form on the joint distribution of the underlying DAG. Instead, it relies on a linear operator on reproducing kernel Hilbert spaces to evaluate conditional independence. However, a fully nonparametric approach would involve conditioning on a large number of random variables, subjecting it to the curse of dimensionality. To solve this problem, nonlinear sufficient dimension reduction is applied to reduce the number of variables before evaluating the conditional independence. An estimator is developed for the DAG, based on a linear operator that characterizes conditional independence, and the consistency and convergence rates of this estimator are established, as well as the uniform consistency of the estimated Markov equivalence class. A modified PC-algorithm is introduced to implement the estimation procedure efficiently, such that the complexity depends on the sparseness of the underlying true DAG. The effectiveness of the methodology is demonstrated through simulations and a real data analysis.

C1199: **Belted and ensembled neural network for linear and nonlinear sufficient dimension reduction**

Presenter: **Yin Tang**, University of Kentucky, United States

Co-authors: Bing Li

The aim is to introduce a unified, flexible, and easy-to-implement framework of sufficient dimension reduction that can accommodate both linear and nonlinear dimension reduction, and both the conditional distribution and the conditional mean as the targets of estimation. This unified framework is achieved by a specially structured neural network – the belted and ensembled neural network (BENN) – that consists of a narrow latent layer, which is called the belt, and a family of transformations of the response, which is called the ensemble. By strategically placing the belt at different layers of the neural network, linear or nonlinear sufficient dimension reduction is achieved, and by choosing the appropriate transformation families, dimension reduction is achieved for the conditional distribution or the conditional mean. Moreover, thanks to the advantage of the neural network, the method is very fast to compute, overcoming a computation bottleneck of the traditional sufficient dimension reduction estimators, which involves the inversion of a matrix of dimension either p or n . The algorithm and convergence rate of the method are developed, compared with existing sufficient dimension reduction methods, and applied to two data examples.

C1215: **Order determination in sufficient dimension reduction for multivariate time series**

Presenter: **Andreas Artemiou**, University of Limassol, Cyprus

Co-authors: Amal Alqarni

The focus is on robust sufficient dimension reduction in multivariate time series. The first effort to determine the dimension of the dimension reduction subspace being estimated using an information criteria approach is discussed.

C1278: **On non-redundant and linear operator-based nonlinear dimension reduction**

Presenter: **Wei Luo**, Zhejiang University, China

Co-authors: Zhoufu Ye

Kernel principal component analysis (KPCA), a popular nonlinear dimension reduction technique, has the redundancy issue that each kernel principal component can be a measurable function of the preceding components. This harms the effectiveness of dimension reduction and leaves the dimension of the reduced data a heuristic choice. The purpose is to rebuild the theory of nonlinear dimension reduction centered on recovering the sigma-field of the original data, and, using appropriate linear operators between RKHSs, two sequential dimension reduction methods are proposed that address the redundancy issue, maintain the same level of computational complexity as KPCA, and rely on more plausible assumptions regarding the singularity of the original data. Compared with the existing nonlinear dimension reduction methods that also address the redundancy issue, the methods enjoy the parametric asymptotic rate and do not specify distributions on the reduced data, thereby preserving other patterns, if any, of the original data. By constructing a measure of the exhaustiveness of the reduced data, consistent order determination is also provided for these methods. Some numerical studies are presented at the end. A novel characterization of conditional mean independence is involved, which may attract independent research interest.

CO291 Room BCB 402 STATISTICAL METHODS FOR LARGE-SCALE BIOMEDICAL DATA ANALYSIS

Chair: Shuting Shen

C1211: **All-in-one rare-variant analysis tool for biobank-scale whole-genome sequencing data**

Presenter: **Zilin Li**, Northeast Normal University, China

Large-scale whole-genome sequencing (WGS) studies have enabled the analysis of rare variant associations with complex human diseases and

traits. Variant set analysis is a powerful approach to studying rare variant associations. However, existing methods have a limited ability to define the variant set in the genome, especially for the noncoding genome. A computationally efficient and robust genetic variant association-detection framework, STAARpipeline, is proposed to automatically annotate a WGS study and perform flexible rare variant association analysis, including functional category-based gene-centric analysis and fixed-window and dynamic-window-based non-gene-centric analysis by incorporating variant functional annotations. STAARpipeline also provides analytical follow-up of dissecting association signals independent of known variants via conditional analysis. The STAARpipeline is applied to analyze Alzheimer's disease (AD) in 459,216 samples from the UK Biobank. All analyses scale well in computation time and memory. Several potentially new significant associations with AD are discovered. In summary, STAARpipeline is a powerful and resource-efficient tool for association analysis of biobank-scale WGS studies.

C1047: Generalized heterogeneous functional model with applications to large-scale mobile health data

Presenter: **Fei Xue**, Purdue University, United States

Co-authors: Xiaojing Sun, Bingxin Zhao

With the increasing availability of large-scale mobile health data, strong associations have been found between physical activity and various diseases. However, accurately capturing this complex relationship is challenging, possibly because it varies across different subgroups of subjects, especially in large-scale datasets. To fill this gap, a generalized heterogeneous functional method is proposed, which simultaneously estimates functional effects and identifies subgroups within the generalized functional regression framework. The proposed method captures subgroup-specific functional relationships between physical activity and diseases, providing a more nuanced understanding of these associations. Additionally, a pre-clustering method that enhances computational efficiency for large-scale data through a finer partition of subjects compared to true subgroups is introduced. A testing procedure is further developed to assess whether the identified subgroups exhibit truly distinct functional effects and whether heterogeneity exists across the entire population. In the real data application, the impact of physical activity is examined on the risk of mental disorders and Parkinson's disease using the UK Biobank dataset, which includes over 79,000 participants. The proposed method outperforms existing methods in future-day prediction accuracy, identifying four subgroups for mental disorder outcomes and three subgroups for Parkinson's disease diagnosis.

C1348: A new Mendelian randomization method integrated with AlphaFold3 for 3D structure prediction

Presenter: **Zhonghua Liu**, Columbia University, United States

Hidden confounding biases hinder the identification of causal protein biomarkers for Alzheimer's disease in non-randomized studies. While Mendelian randomization (MR) can mitigate these biases using protein quantitative trait loci (pQTLs) as instrumental variables, some pQTLs violate core assumptions, leading to biased conclusions. To address this, MR-SPI is proposed, a novel MR method that selects valid pQTL instruments using Leo Tolstoy's Anna Karenina principle and performs robust post-selection inference. Integrating MR-SPI with AlphaFold3, a computational pipeline is developed to identify causal protein biomarkers and predict 3D structural changes. Applied to genome-wide proteomics data from 54,306 UK Biobank participants and 455,258 subjects (71,880 cases and 383,378 controls) for a genome-wide association study of Alzheimer's disease, seven proteins are identified (TREM2, PILRB, PILRA, EPHA1, CD33, RET, and CD55) with structural alterations due to missense mutations. These findings offer insights into the etiology and potential drug targets for Alzheimer's disease.

C1430: MiLC to account for unobserved confounding and reduce false discoveries in microbiome research

Presenter: **Siyuan Ma**, Vanderbilt University Medical Center, United States

Recent research has highlighted false discoveries in microbiome studies, particularly in differential abundance (DA) analyses. While data compositionality has received attention, it is demonstrated that unobserved confounding (e.g., population heterogeneity, recent antibiotic use, or seasonal dietary changes) can be an even stronger driver of false discoveries. Using real-data evidence, it is shown that unobserved confounding inflates false discoveries in microbiome DA more than data compositionality. To address this, a novel factor-modeling regression method is introduced, Microbiome Latent Confounder DA (MiLC), to estimate unobserved confounding factors and control false discoveries. MiLC can be applied to both relative abundance and read count microbiome data. Its performance is validated in controlling false discoveries, relative to existing methods, using extensive simulation- and real-data-based benchmarking. Results highlight the critical need to correct for hidden confounders, offering a more reliable framework for microbiome DA analyses and ultimately improving the robustness of microbiome research findings.

CO105 Room BCB 403 COMPLEX PROBLEMS IN CAUSAL INFERENCE

Chair: Michael Daniels

C0423: Weighting methods for survivor average causal effect estimation in cluster-randomized trials

Presenter: **Nandita Mitra**, University of Pennsylvania, United States

Co-authors: Dane Isenberg, Fan Li, Michael Harhay

Patient-centered outcomes, such as quality of life and length of hospital stay, are often the focus of clinical studies. However, elderly or critically ill clinical trial participants may have truncated or undefined non-mortality outcomes if they do not survive through the measurement time point. To address truncation by death, the survivor average causal effect (SACE) has been proposed as a causally interpretable subgroup treatment effect defined under the principal stratification framework. However, most methods for estimating SACE have been developed in the context of individually-randomized trials. Only limited discussions have centered on cluster-randomized trials (CRTs), where methods typically involve strong distributional assumptions for outcome modeling. Two weighting methods are proposed to estimate SACE in CRTs that obviate the need for potentially complicated outcome distribution modeling. Assumptions that address latent clustering effects to enable point identification of SACE are established, and computationally-efficient asymptotic variance estimators are provided for each weighting estimator. In simulations, the weighting estimators are evaluated, demonstrating their finite-sample operating characteristics and robustness to certain departures from the identification assumptions. The methods are illustrated using data from a CRT to assess the impact of a sedation protocol on mechanical ventilation among children with acute respiratory failure.

C0489: A Bayesian nonparametric approach to mediation and spillover effect with multiple mediators in cluster-randomized trials

Presenter: **Fan Li**, Yale University, United States

Co-authors: Yuki Ohnishi

Cluster randomized trials (CRTs) with multiple unstructured mediators present significant methodological challenges for causal inference due to within-cluster correlation, interference among units, and the complexity introduced by multiple mediators. Existing causal mediation methods often fall short in simultaneously addressing these complexities, particularly in disentangling mediator-specific effects under interference that are central to studying complex mechanisms. To address this gap, new causal estimands are proposed for spillover mediation effects that differentiate the roles of each individual's own mediator and the spillover effects resulting from interactions among individuals within the same cluster. Identification results are established for each estimand and, to flexibly model the complex data structures inherent in CRTs, a new Bayesian nonparametric prior is developed - the nested dependent Dirichlet process mixture - designed to flexibly capture the outcome and mediator surfaces at different levels. Extensive simulations are conducted across various scenarios to evaluate the frequentist performance of the methods, compare them with a Bayesian parametric counterpart, and illustrate the new methods in an analysis of a completed CRT.

C0468: A case study of causal mediation using Bayesian nonparametrics and semiparametric corrections

Presenter: **Yuhua Zhang**, University of Florida, United States

A Bayesian nonparametric approach is proposed using a truncated enriched Dirichlet process mixture (EDPM) model to estimate natural direct (NDE) and indirect (NIE) effects in causal mediation analyses in the presence of post-treatment confounders. An efficient cluster reallocation

Metropolis-Hasting algorithm is introduced to improve mixing in the blocked Gibbs sampler. A one-step posterior correction is implemented based on the efficient influence function for the setting. This post-processing step solves a critical problem in Bayesian nonparametrics: How to obtain reliable estimates and posteriors for a specific causal estimand of interest (the NDE and NIE) with excellent frequentist properties, such as correct coverage, from a model designed for complex joint distributions. Simulation studies are conducted to assess the method's performance, and it is applied to evaluate causal mediation effects in a weight management clinical trial.

C0521: Conformal inference for multivariate mixed outcomes via fairness-constrained optimal transport

Presenter: **Larry Han**, Northeastern University, United States

Co-authors: Chenyin Gao

Many decisions require reasoning over multiple outcomes simultaneously, especially when outcomes are of mixed type (e.g., continuous and discrete). Existing conformal inference methods for multivariate prediction often assume continuous outcomes, overlook fairness across subgroups, or yield prediction regions that are difficult to interpret. A conditional latent highest density region (CL-HDR) is introduced, a new framework for conformal prediction with multivariate mixed outcomes that ensures finite-sample coverage, supports customizable fairness constraints, and yields efficient prediction regions with minimal volume. The method leverages optimal transport and normalizing flows to construct multivariate ranks and uses input convex neural networks to approximate the transport map. To guarantee group-conditional coverage, a functional synchronization procedure is proposed based on Wasserstein barycenters, enabling fairness objectives to be encoded directly into the prediction set construction. Across synthetic and real-world datasets, CL-HDR produces smaller, more interpretable prediction regions than existing approaches while achieving subgroup-level validity.

C0248 Room BCB 405 RECENT DEVELOPMENTS IN BAYESIAN CLUSTERING

Chair: Alessandro Casa

C0226: Bayesian level set clustering

Presenter: **Miheer Dewaskar**, University of New Mexico, United States

Co-authors: David Buch, David Dunson

Classically, Bayesian clustering interprets each component of a mixture model as a cluster. The inferred clustering posterior is highly sensitive to any inaccuracies in the kernel within each component. As this kernel is made more flexible, problems arise in identifying the underlying clusters in the data. To address this pitfall, a fundamentally different approach is proposed to Bayesian clustering that decouples the problems of clustering and flexible modeling of the data density f . Starting with an arbitrary Bayesian model for f and a loss function for defining clusters based on f , a Bayesian decision-theoretic framework is developed for density-based clustering. Within this framework, a Bayesian level set clustering method is developed to cluster data into connected components of a level set of f . Theoretical support is provided, including clustering consistency, and performance is highlighted in a variety of simulated examples. An application to astronomical data illustrates improvements over the popular DBSCAN algorithm in terms of accuracy, insensitivity to tuning parameters, and providing uncertainty quantification.

C0472: Bayesian multi-resolution clustering via infinite latent factors

Presenter: **Lorenzo Schiavon**, Ca Foscari University of Venice, Italy

Co-authors: Mattia Stival

In many scientific fields, a critical task is to cluster subjects based on a potentially vast set of features. A fundamental challenge in model-based clustering is the trade-off between the resolution of the inferred clusters and the parsimony of the model. Current nonparametric approaches often require a pre-specified resolution level, demanding extensive parameterization to capture fine-grained structures and offering no mechanism to explore cluster hierarchies. To overcome these limitations, a novel multi-resolution clustering approach is introduced using an infinite mixture model with kernels organized in a multiscale framework. The method, through a careful specification of mixture weights, naturally incorporates exogenous information to guide the formation of cluster hierarchies while maintaining flexibility. The theoretical properties of the model are investigated, and an elegant and parsimonious formulation is proposed based on an infinite factorization, which allows for efficient posterior inference via a Gibbs sampler. The practical advantages of the approach are shown on synthetic data and through challenging real-world applications, revealing multi-level grouping patterns in survey responses and gene expression data.

C0476: Informed partition models for dependent random partitions

Presenter: **Sally Paganin**, The Ohio State University, United States

Co-authors: Garritt Page, Fernando Quintana

Model-based clustering is a powerful tool often used to discover hidden structure in data by grouping observational units that exhibit similar response values. Recently, clustering methods have been developed that allow the inclusion of an initial partition of the data informed by expert opinions, starting from a probability distribution on the space of partitions. Then, using some similarity criteria, partitions different from the initial one are down-weighted, i.e., they are assigned reduced probabilities. A different perspective is taken, and the probability that each unit follows the initial partition via auxiliary variables is modeled. The informed partition model provides flexibility to include varying levels of uncertainty to any subset of the partition (i.e., locally weighted prior information). Additionally, it can accommodate settings with multiple dependent partitions, such as temporal or multi-view data. Theoretical properties of the proposed construction are explored, which can be useful for prior elicitation. The gains in prior specification flexibility are illustrated via simulation studies and an application to a dataset concerning the spatiotemporal evolution of PM10 measurements in Germany.

C0676: Learning pathways of life events: A sequential allocation Bayesian model with label tracking

Presenter: **Beatrice Franzolini**, Bocconi University, Italy

Co-authors: Andrea Cremaschi, Raffaella Piccarreta

The analysis of socio-economic data that tracks life events, such as marriage, childbirth, and employment, can deepen the understanding of representative life trajectories in contemporary societies. However, this typically requires clustering methods for multivariate longitudinal categorical data, and existing model-based methods in the statistical literature remain rather limited in this regard. A flexible Bayesian nonparametric framework is introduced for modeling multiple dependent categorical variables observed over time through evolving latent life stages. The primary goal is to identify the stages individuals pass through during life and to characterize distinct dynamic behavioral patterns. A key methodological innovation is that each life stage is defined by time-invariant parameters, ensuring a consistent interpretation across all time points. Unlike modern Bayesian dynamic clustering approaches, which re-estimate cluster characteristics at every time step, the strategy greatly improves computational efficiency and enables meaningful longitudinal comparisons. In addition, the model captures temporal dynamics via a time-varying partition of the study population that incorporates both abrupt structural change points and individual-level transitions between life stages. The methodology is illustrated using longitudinal data from the Italian Institute of Statistics, demonstrating its ability to reveal both individual behaviors and broader societal shifts.

C1038: Understanding uncertainty in Bayesian cluster analysis

Presenter: **Cecilia Balocchi**, University of Edinburgh, United Kingdom

Co-authors: Sara Wade

The Bayesian approach to clustering is often appreciated for its ability to provide uncertainty in the partition structure. However, summarizing the posterior distribution over the clustering structure can be challenging, due to the discrete, unordered nature and massive dimension of the space.

While recent advancements provide a single clustering estimate to represent the posterior, this ignores uncertainty and may even be unrepresentative in instances where the posterior is multimodal. To enhance the understanding of uncertainty, a Wasserstein approximation for Bayesian clustering (WASABI) is proposed, which summarizes the posterior samples with not one, but multiple clustering estimates, each corresponding to a different part of the space of partitions that receives substantial posterior mass. Specifically, such clustering estimates are found by approximating the posterior distribution in a Wasserstein distance sense, equipped with a suitable metric on the partition space. An interesting byproduct is that a locally optimal solution to this problem can be found using a k-medoids-like algorithm on the partition space to divide the posterior samples into different groups, each represented by one of the clustering estimates. Using both synthetic and real datasets, it is shown that the proposal helps to improve the understanding of uncertainty, particularly when the data clusters are not well separated or when the employed model is misspecified.

CO104 Room BCB 406 RECENT ADVANCEMENTS IN STATISTICAL NETWORK ANALYSIS
Chair: Jonathan Stewart
C0353: Minorization-maximization-based estimation for network models with parameter vectors of increasing dimension

Presenter: **Cornelius Fritz**, Trinity College Dublin, Ireland

Co-authors: Michael Schweinberger, David Hunter

Large and complex network data necessitate complex models accommodating local and global dependence. Local dependence refers, e.g., to transitive clustering based on common partners, implying the knowledge about other population members' connections. On the other hand, global dependence governs the general propensity to interact with other population members regardless of sharing a common neighborhood. A general framework is introduced to capture both types of dependencies that give rise to high-dimensional models for network data. Standard algorithms to tackle such computational problems are based on Metropolis-Hastings Monte Carlo or Newton-Raphson Methods, which perform poorly in high-dimensional settings. A minorization-maximization (MM) method is introduced for convex objective functions to alleviate this scaling issue. Quasi-Newton acceleration methods are employed to speed up the convergence of the algorithm. Moreover, a penalty is introduced to bypass issues of unidentifiable coefficients. In several applications, the performance of the algorithms is exhibited in comparison to currently available algorithms.

C0793: Assessing the explanatory power of high-dimensional node-level covariates on network structure

Presenter: **Alexander Fuchs-Kreiss**, Leipzig University, Germany

Co-authors: Keith Levin

The purpose is to consider the problem of selecting from a high-dimensional set of covariates those that are explanatory of an observed network. Specifically, a network comprising vertices and undirected links is observed between them. In addition, for each vertex, a high-dimensional vector of covariates is observed (its dimension can exceed the number of vertices). To assess which of these covariates are correlated with the network structure, it is assumed that the network is generated by a random dot product graph (RDPG), and the aim is to understand if some covariates are related to the latent positions in the RDPG. This is achieved in three ways. Firstly, a model is proposed, and LASSO estimation is used. The unidentifiability in the RDPG translates to a group LASSO penalty. In a second approach, canonical correlation analysis (CCA) is used to quantify the strength of the relation between the latent positions and the covariates. Since the latent positions are unobserved, both methods are applied to the estimated latent positions. Thirdly, as an alternative, CCA is studied between the covariates and the rows of the adjacency matrix directly. This avoids the need for latent position estimation. For all the methods, the convergence of the estimators is rigorously shown. In a simulation study, it is shown that permutation tests based on the methods have good power in detecting relevant covariates from a high-dimensional set of covariates.

C0863: Steins method of moment estimators for local dependency exponential random graph models

Presenter: **Gesine Reinert**, Oxford University, United Kingdom

Co-authors: Adrian Fischer, Wenkai Xu

Providing theoretical guarantees for parameter estimation in exponential random graph models is a largely open problem. While maximum likelihood estimation has theoretical guarantees in principle, verifying the assumptions for these guarantees to hold can be very difficult. Moreover, in complex networks, numerical maximum likelihood estimation is computer-intensive and may not converge in a reasonable time. To ameliorate this issue, local dependency exponential random graph models have been introduced, which assume that the network consists of many independent exponential random graphs. In this setting, progress towards maximum likelihood estimation has been made. However, the estimation is still computer-intensive. Instead, the proposal is to use so-called Stein estimators: The Stein characterizations are used to obtain new estimators for local dependency exponential random graph models.

C0969: Clustering and inference for very sparse diverse multiplex networks

Presenter: **Marianna Pensky**, University of Central Florida, United States

The focus is on the DIVERse MultiPLEx Generalized Random Dot Product Graph (DIMPLE-GRDPG) network model, where all layers of the network have the same collection of nodes and follow the Generalized Random Dot Product Graph (GRDPG) model. In addition, all layers can be partitioned into groups such that the layers in the same group are embedded in the same ambient subspace, but otherwise all matrices of connection probabilities can be different. While this is already a very difficult model, it is also assumed that layers of the network are very sparse. Tensor-based approaches are used to recover the groups of layers in such a network, and subsequently estimate the ambient subspaces.

CO157 Room BCB 407 CAUSAL INFERENCE AND PERSONALIZED MEDICINE
Chair: Chan Park
C1123: The conflict graph design: Estimating causal effects under network interference

Presenter: **Christopher Harshaw**, Columbia University, United States

Co-authors: Vardis Kandiros, Charilaos Pipis, Constantinos Daskalakis

From political science and economics to public health and corporate strategy, the randomized experiment is a widely used methodological tool for estimating causal effects. In the past 15 years or so, there has been a growing interest in network experiments, where subjects are presumed to be interacting in the experiment and their interactions are of substantive interest. While the literature on interference has focused primarily on unbiased and consistent estimation, designing randomized network experiments to ensure tight rates of convergence is relatively under-explored. Not only are the optimal rates of estimation for different causal effects under interference an open question, but previously proposed designs are created in an ad-hoc fashion. A new experimental design is presented for network experiments called the "Conflict Graph Design", which, given a pre-specified causal effect of interest and the underlying network, produces a randomization over treatment assignment with the goal of increasing the precision of effect estimation. Not only does this experiment design attain improved rates of consistency for several causal effects of interest, but it also provides a unifying approach to designing network experiments. Consistent variance estimators and asymptotically valid confidence intervals are also provided, which facilitate inference of the causal effect under investigation.

C1175: Learning robust treatment rules for censored data

Presenter: **Yifan Cui**, Zhejiang University, China

There is a fast-growing literature on estimating optimal treatment rules directly by maximizing the expected outcome. In biomedical studies and operations applications, censored survival outcome is frequently observed, in which case the restricted mean survival time and survival probability are of great interest. Two robust criteria are proposed for learning optimal treatment rules with censored survival outcomes; the former one targets at an optimal treatment rule maximizing the restricted mean survival time, where the restriction is specified by a given quantile such as median; the latter one targets at an optimal treatment rule maximizing buffered survival probabilities, where the predetermined threshold is adjusted to

account the restricted mean survival time. Theoretical justifications are provided for the proposed optimal treatment rules, and a sampling-based difference-of-convex algorithm is developed for learning them. In simulation studies, the estimators show improved performance compared to existing methods. The proposed method is also demonstrated using AIDS clinical trial data.

C1270: Resource-efficient policy targeting under heterogeneous partial interference

Presenter: **Elena Dal Torrione**, Yale University, United States

Co-authors: Chan Park, Laura Forastiere

In many empirical studies, units are interconnected, and a unit's outcome may depend on the treatment of others, leading to interference. When interference is heterogeneous, treating individuals with specific characteristics can influence the population average outcome differently, either through their direct response or their impact on others. For instance, policymakers may minimize resource use by vaccinating individuals identified as superspreaders to achieve a target reduction in disease incidence. Under heterogeneous clustered interference, we propose a method to estimate optimal stochastic treatment allocations, in which an individual's treatment probability is determined by both individual- and cluster-level covariates. The approach minimizes the expected marginal treatment probability within a cluster while ensuring a specified outcome level is met. To evaluate the methodology, theoretical guarantees are provided, analyzing how the excess risk bound depends on function class complexity and cluster size. Additionally, a simulation study is conducted, and the method is applied to a water, sanitation, and hygiene (WASH) intervention in Senegal. The estimated policy is compared to alternative approaches, expecting the method to achieve greater resource efficiency compared to policies with homogeneous treatment probabilities within clusters.

C1448: Inference on local variable importance measures for heterogeneous treatment effects

Presenter: **Pawel Morzywolek**, University of Copenhagen, Denmark

Co-authors: Peter Gilbert, Alex Luedtke

An inferential framework is provided to assess variable importance for heterogeneous treatment effects. This assessment is especially useful in high-risk domains such as medicine, where decision makers hesitate to rely on black-box treatment recommendation algorithms. The variable importance measures we consider are local in that they may differ across individuals, while the inference is global in that it tests whether a given variable is important for any individual. The approach builds on recent developments in semiparametric theory for function-valued parameters and is valid even when statistical machine learning algorithms are employed to quantify treatment effect heterogeneity.

CO303 Room BCB 408 ADVANCES IN BAYESIAN METHODS AND COMPUTATIONS

Chair: Maria De Iorio

C0656: MCMC inference for latent semi-Markov point processes

Presenter: **Rosario Barone**, Università Cattolica del Sacro Cuore, Italy

This contribution develops a Bayesian framework for inference on continuous-time point processes driven by latent semi-Markov dynamics. Unlike standard Markov-modulated models, the latent process is allowed to exhibit non-exponential sojourn times, thus capturing duration dependence and temporal heterogeneity. The framework supports a wide range of conditional intensity functions, including homogeneous Poisson, covariate-modulated, self-exciting, and self-correcting forms. Posterior inference is performed using a Metropolis-within-Gibbs sampler that combines a forward-backward algorithm tailored for semi-Markov models with a uniformization-based path-sampling strategy. This approach enables exact inference under the continuous-time model without requiring discretization or approximation steps. A simulation study confirms the method's ability to recover both latent trajectories and model parameters with high accuracy, while preserving computational efficiency. An empirical application to drug-related arrests in Chicago demonstrates that the semi-Markov specification improves the detection of latent regime shifts compared to Markovian benchmarks, particularly under self-exciting dynamics. The proposed approach offers a flexible and efficient solution for modeling event-driven data with irregular timing and complex latent structure.

C1045: Bayesian mixture models with repulsive and attractive atoms

Presenter: **Alessandra Guglielmi**, Politecnico di Milano, Italy

The study of almost surely discrete random probability measures is an active line of research in Bayesian non-parametrics. The idea of assuming interaction across the atoms of the random probability measure has recently spurred significant interest in the context of Bayesian mixture models. This allows the definition of priors that encourage well-separated and interpretable clusters. A unified framework is provided for the construction and the Bayesian analysis of random probability measures with interacting atoms, encompassing both repulsive and attractive behaviours. Specifically, closed-form expressions are derived for the posterior distribution, the marginal and predictive distributions, previously unavailable except for the case of measures with i.i.d. atoms. It is shown how these quantities are fundamental for both prior elicitation and developing new posterior simulation algorithms for hierarchical mixture models. Results are obtained without any assumption on the finite point process governing the atoms of the random measure. The treatment is specialized to the classes of Poisson, Gibbs, and determinantal point processes, as well as in the case of shot-noise Cox processes. Finally, the modelling strategies are seen on simulated and real datasets.

C0811: Filtering Wright-Fisher diffusions via discrete dual processes

Presenter: **Guillaume Kon Kam King**, Université Paris-Saclay, INRAE, France

Co-authors: Paul Jenkins, Matteo Ruggiero

Exact inference for hidden Markov models requires the evaluation of all distributions of interest, filtering, prediction, smoothing, and likelihood with a finite computational effort. We provide sufficient conditions for exact inference for a class of hidden Markov models on general state spaces, given a set of discretely collected indirect observations nonlinearly linked to the signal, and a set of practical inference algorithms. The conditions we obtain are concerned with the existence of a certain type of dual process, which is an auxiliary process embedded in the time reversal of the signal, that in turn allows us to represent the distributions and functions of interest as countable mixtures of elementary densities or products thereof. We explore the applicability of this strategy for exact inference on Wright-Fisher diffusions, with or without selection.

C1210: Modeling pairwise comparison data with cyclic and acyclic structures under a Bayesian framework

Presenter: **Hisaya Okahara**, Tokyo University of Science, Japan

Co-authors: Tomoyuki Nakagawa, Shonosuke Sugawara

Pairwise comparison data frequently arises in diverse applications such as sports analytics, preference learning, and human feedback evaluation. Classical models, such as the Bradley-Terry model, assume that preferences can be explained by latent scores associated with each item, leading to a globally consistent ranking. However, real-world data often exhibit cyclic or intransitive patterns that cannot be captured by score-based approaches alone. A statistical modeling framework that extends classical pairwise comparison models is introduced to flexibly accommodate such complex patterns while preserving interpretability. The approach is formulated under a Bayesian logistic regression setting, where a Gibbs sampling framework enables efficient posterior computation. Finally, through numerical experiments, the proposed framework is illustrated to be capable of capturing cyclic structures observed in practice.

CO202 Room BCB 409 DEVELOPMENTS IN SPATIO-TEMPORAL DISEASE MAPPING AND SURVEILLANCE

Chair: Andrew Lawson

C0188: Ensemble Kalman filtering for Bayesian disease map surveillance

Presenter: **Andrew Lawson**, Medical University of South Carolina, United States

Kalman filtering has been commonly used for time series modeling and the prediction of time-based systems. In the exploration of application

to disease mapping, the conversion of count data is considered from small areas to a Gaussian-like variable: $\text{Log}(\text{SMR})$. Because the $\text{log}(\text{SMR})$ behaves closely like a Gaussian, the use of KF methods is explored for space-time data. However, for fitting purposes, especially for online surveillance, the need to invert large matrices is unattractive. Instead, the use of ensemble KF is explored, which bypasses the inversion by using simulation-based updating, closely related to particle filtering. It has been found that very fast updating of system and measurement models can be made without (much!) recourse to MCMC, and so it would appear to provide a useful tool in applications where fast updating is required, such as during large infectious disease outbreaks.

C0247: **Spatial disaggregation using a Bayesian low-rank geoadditive model**

Presenter: **Christel Faes**, Hasselt University, Belgium

A novel Bayesian spatial disaggregation model is presented for count data that integrates smooth covariate effects and spatial dependencies within a unified framework. The model leverages penalized splines to flexibly capture nonlinear relationships with covariates. For modeling spatial correlation, a spline-based low-rank kriging approximation is adopted, which improves computational tractability. Additionally, Laplace approximation is used for posterior inference, offering substantial efficiency gains over traditional Markov chain Monte Carlo (MCMC) techniques. Two estimation strategies are explored: One based on the exact likelihood and another using a spatially discrete approximation to further enhance computational performance. Through simulation studies, it is shown that both approaches yield accurate estimates, with the approximate method delivering notable computational savings. The model is applied to disaggregate disease incidence data, generating high-resolution risk maps. Overall, the approach provides a flexible, scalable, and practical framework for spatial disaggregation in geostatistical applications involving count data.

C0460: **A Bayesian spatial model for survey-based ordinal data**

Presenter: **Ana Corberan-Vallet**, University of Valencia, Spain

Co-authors: Miguel Angel Beltran-Sanchez, Miguel Angel Martinez-Beneito

Health surveys allow exploring health indicators that are of great value from a public health point of view. These indicators are usually coded as ordinal variables and depend on covariates associated with individuals. A Bayesian individual-level model is proposed for small-area estimation of survey-based health indicators. A categorical likelihood is used to describe the ordinal data. At the second level of the model hierarchy, the cumulative probabilities of the different categories are modeled, taking into account possible covariate effects as well as spatial dependence among areas. Post-stratification of the results allows extrapolating the results to any administrative area division, even for small areas. Finally, a multivariate extension of the model is presented that allows for the joint study of the sets of response variables that are likely to be correlated.

C0746: **Sequential Bayesian spatiotemporal outbreak detection**

Presenter: **Frank Zou**, University of Central Florida, United States

Early online outbreak detection for an epidemic is vital for disease-control authorities to make policies for the protection of public health and normal socioeconomic functions. Modern public health streaming surveillance data are often collected from multiple data sources, exhibiting spatio-temporal interdependence and imbalance issues. To address those issues, a Bayesian online spatiotemporal outbreak detection is proposed with prior updating and p-value adaptation (BOSTON-PUPA) procedure. Using sequential p-value combinations, this iterative procedure involves the generalized Poisson distribution (GPD) model and supports synchronous surveillance over multiple locations, with a controlled false detection rate as well as high sensitivity against outbreaks in a wide range of signal-to-noise ratios. In the simulation study, several popular combined p-value methods in the BOSTON-PUPA procedure are employed and compared based on sensitivity, specificity, false detection rate, and delay before making recommendations. The method is illustrated by detecting the outbreaks in the real COVID-19 daily case count data in Massachusetts counties in 2020.

CC444 Room BCB 311 SHORT TALKS: CMSTATISTICS I

Chair: Silvia Montagna

C1507: **Robust semiparametric inference for Bayesian additive regression trees**

Presenter: **Ruixuan Liu**, Chinese University of Hong Kong, Hong Kong

A semiparametric framework is developed for inference on the mean response in missing-data settings, using a corrected posterior distribution. Our approach is tailored to Bayesian Additive Regression Trees (BART), a powerful predictive method, but its nonsmoothness complicates asymptotic theory for multidimensional covariates. When using BART combined with Bayesian bootstrap weights, we establish a new Bernstein von Mises theorem and show that the limit distribution generally contains a bias term. To address this, we introduce RoBART, a posterior bias-correction that robustifies BART for valid inference on the mean response. Monte Carlo studies support our theory, demonstrating reduced bias and improved coverage relative to existing procedures using BART.

C1509: **A kernel-based test for the proportional hazards assumption**

Presenter: **Tamara Fernandez**, Universidad Adolfo Ibanez, Chile

Co-authors: Nicolas Rivera, Merle Munko

A novel kernel-based test is presented for assessing the proportional hazards assumption. The proposed test evaluates whether deviations from the Cox model, captured through functions in a reproducing kernel Hilbert space (RKHS), lead to a better approximation of the underlying hazard. Unlike standard kernel tests, our approach faces an additional challenge: it requires estimating the Cox model parameters under the null hypothesis as part of the testing procedure. This estimation step implies that conventional resampling strategies, such as Wild Bootstrap, are no longer valid. We outline the derivation of the test statistic and introduce a corrected wild bootstrap method that accounts for parameter estimation and provides valid inference.

C1524: **Mixture of state space models for compositional data with an application to urban mobility analysis**

Presenter: **Andrea Panarotto**, Department of Statistical Sciences, University of Padova, Italy

Co-authors: Manuela Cattelan, Ruggero Bellio

Capturing the dependency between successive compositional observations requires models that account for the constrained nature of the data. A framework is introduced for the analysis and clustering of compositional time series, where the trajectories evolve within the simplex. The approach integrates a state-space representation of compositional dynamics into a model-based clustering framework, enabling the identification of groups of trajectories sharing similar temporal patterns. The model can be extended through a mixture of experts model, allowing trajectory-level covariates to influence the component weights. The methodology is motivated by an application to urban mobility, where individuals' movements are represented in the simplex by the proportions of types of roads in their surroundings. This formulation provides a data-driven way to aggregate individual travel behaviors into population-wide mobility patterns, offering new insights into how people interact with different urban environments over time.

C1523: **Efficient estimation with accelerated gap-time model for recurrent events**

Presenter: **Akim Adekpedjou**, Missouri University of Science and Technology, United States

Co-authors: Emmanuel Djegou

The accelerated failure time model is a regression-like model that relates, linearly, a set of covariates to the logarithm of failure times with right-censored data. The accelerated gap-time (AGT) model relates, in the same manner, a set of gap-times to covariates, albeit with recurrent events. Both models have the ability to accelerate or decelerate event occurrences via the covariates and provide a quite direct physical interpretation of

event time. They are also a good alternative to the popular Cox model. They have broad applications in engineering and biomedical studies. In this talk, we present the AGT model for recurrent events under a family of effective age models that captures various aging patterns. Estimating parameters with such models for any of these data types is quite challenging due to the non-monotonicity of the score function. To alleviate this, we propose a weighted efficient score function motivated by a family of parametric baseline hazard sub-models. This approach yields estimators with nice large-sample properties. We present results of the procedures for various effective age processes. Simulation studies show good approximations to the true. Limitations and extensions are presented. The procedures are illustrated with a biomedical dataset.

Saturday 13.12.2025

16:50 - 18:55

Parallel Session E – CFE-CMStatistics 2025

CI246 Room BCB 307 CSDA STATISTICAL DATA SCIENCE**Chair: Cristian Gatu****C0277: Training a classifier via semi-supervised learning****Presenter:** Geoffrey McLachlan, University of Queensland, Australia

There has been much increasing attention to semi-supervised learning (SSL) approaches in machine learning for forming a classifier in situations where the training data for a classifier consists of a limited number of classified observations but a much larger number of unclassified observations whose labels denoting their class of origin are unknown. The surprising result of a prior study is considered further, that a classifier formed from a partially classified sample can actually have a smaller expected error rate than if the sample were completely classified. This rather paradoxical outcome is able to be achieved by introducing a framework with a missingness mechanism for the missing labels of the unclassified observations. Within this framework, the conditional probability $q(y)$ that an observation with feature vector y has a missing label is taken to be a logistic model with covariate equal to an entropy-based measurement $e(y)$. The extension of the model is considered for $q(y)$ to the two-component mixture model, $c + (1-c)q(y)$, where c is the probability that a feature has a label that is missing completely at random (MCAR). The asymptotic relative efficiency of the estimated Bayes' classifier is derived. Results are presented to show how its relative efficiency falls away as c increases. The focus is on two classes in which y has a multivariate normal distribution.

C0455: Robust forecasting with machine learning**Presenter:** Christophe Croux, KU Leuven, Belgium

The impact of outliers on time series prediction is discussed. The use of ARMA models for robust estimation and prediction is well studied. For nonlinear models, however, machine learning methods are a suitable alternative. Popular machine learning methods such as random forests, XGBoost, and LightGBM can also be used in a time series context. Their predictive performance is examined in the presence of outliers. Moreover, it is investigated how these methods can be made more robust by changing the loss function and adding a filtering step.

C1350: On MCMC mixing for predictive inference under unidentified transformation models**Presenter:** Catherine Liu, The Hong Kong Polytechnic University, Hong Kong

Reliable Bayesian predictive inference has long been an open problem under unidentified transformation models, since the Markov chain Monte Carlo (MCMC) chains of posterior predictive distribution (PPD) values are generally poorly mixed. The aim is to address the poorly mixed PPD value chains under unidentified transformation models through an adaptive scheme for prior adjustment. Specifically, a conception of sufficient informativeness is originated, which explicitly quantifies the information level provided by nonparametric priors, and assesses MCMC mixing by comparison with the within-chain MCMC variance. The prior information level is formulated by a set of hyperparameters induced from the nonparametric prior elicitation with an analytic expression, which is guaranteed by asymptotic theory for the posterior variance under unidentified transformation models. The analytic prior information level consequently drives a hyperparameter tuning procedure to achieve MCMC mixing. The proposed method is general enough to cover various data domains through a multiplicative error working model. Comprehensive simulations and real-world data analysis demonstrate that the method successfully achieves MCMC mixing and outperforms state-of-the-art competitors in predictive capability.

C1477: An algorithmic procedure for solving the generalized minimum information checkerboard copula problem**Presenter:** Ivan Kojadinovic, CNRS UMR 5142 LMA University of Pau, France**Co-authors:** Tommaso Martini

The minimum information copula principle (see Meeuwissen and Bedford, 1997) is a maximum entropy-like approach for finding the least informative copula that satisfies a certain number of expectation constraints specified either from expert knowledge or the available limited data. In this presentation, we first propose a generalization of this principle allowing the inclusion of additional constraints fixing certain higher-order margins of the copula. We next prove that the associated optimization problem has a unique solution under a natural condition. As the latter problem is intractable in general, following the existing literature, we consider its version with all the probability measures involved in its formulation replaced by checkerboard approximations. This amounts to attempting to solve a so-called discrete 1-projection linear problem. We then use the seminal results of Csiszar (1975) to derive an IPFP-like procedure for solving the latter and provide theoretical guarantees for its convergence. We conclude the presentation with numerical experiments in dimensions up to four with substantially finer discretizations than those encountered in the literature.

CO168 Room BCB G07 ADVANCES IN STATISTICAL MACHINE LEARNING FOR MODERN DATA CHALLENGES**Chair: Xiwei Tang****C0372: Online estimation and inference for robust policy evaluation in reinforcement learning****Presenter:** Yichen Zhang, Purdue University, United States

Reinforcement learning has emerged as one of the prominent topics attracting attention in modern statistical learning, with policy evaluation being a key component. Unlike the traditional machine learning literature on this topic, statistical inference is emphasized for the model parameters and value functions of reinforcement learning algorithms. While most existing analyses assume random rewards to follow standard distributions, the concept of robust statistics is embraced in reinforcement learning by simultaneously addressing issues of outlier contamination and heavy-tailed rewards within a unified framework. A fully online robust policy evaluation procedure is developed, and the Bahadur-type representation of the estimator is established. Furthermore, an online procedure is developed to efficiently conduct statistical inference based on the asymptotic distribution. Robust statistics and statistical inference in reinforcement learning are connected, offering a more versatile and reliable approach to online policy evaluation. Finally, the efficacy of the algorithm is validated through numerical experiments conducted in simulations and real-world reinforcement learning experiments.

C1014: Distributional instrumental variable method**Presenter:** Xinwei Shen, University of Washington, United States

The instrumental variable (IV) approach is commonly used to infer causal effects in the presence of unmeasured confounding. Existing methods typically aim to estimate the mean causal effects, whereas a few other methods focus on quantile treatment effects. The aim is to estimate the entire interventional distribution, which yields the classical causal estimands as functionals. A method called distributional instrumental variable (DIV) is proposed, which uses generative modelling in a nonlinear IV setting. Identifiability of the interventional distribution is established under general assumptions, and an 'under-identified' case is demonstrated, where DIV can identify the causal effects while two-step least squares fails to. The empirical results show that the DIV method performs well for a broad range of simulated data, exhibiting advantages over existing IV approaches in terms of the identifiability and estimation error of the mean or quantile treatment effects. Furthermore, DIV is applied to an economic data set to examine the causal relation between institutional quality and economic development, and the results align well with the original study. DIV is also applied to a single-cell data set, where we study the generalizability and stability in predicting gene expression under unseen interventions.

C1398: Heterogeneous transfer learning for high dimensional regression with feature mismatch**Presenter:** Subhadeep Paul, The Ohio State University, United States**Co-authors:** Jae Ho Chang, Massimiliano Russo

The focus is on the problem of transferring knowledge from a source, or proxy, domain to a new target domain for learning a high-dimensional

regression model with possibly different features. Recently, the statistical properties of homogeneous transfer learning have been investigated. However, most homogeneous transfer and multi-task learning methods assume that the target and proxy domains have the same feature space, limiting their practical applicability. In applications, target and proxy feature spaces are frequently inherently different, for example, due to the inability to measure some variables in the target data-poor environments. Conversely, existing heterogeneous transfer learning methods do not provide statistical error guarantees, limiting their utility for scientific discovery. A two-stage method is proposed that involves learning the relationship between the missing and observed features through a projection step in the proxy data and then solving a joint penalized regression optimization problem in the target data. An upper bound is developed on the method's parameter estimation risk and prediction risk, assuming that the proxy and the target domain parameters are sparsely different. The results elucidate how estimation and prediction error depend on the complexity of the model, sample size, the extent of overlap, and correlation between matched and mismatched features.

C1447: Blessing from human-AI interaction: Super reinforcement learning in confounded environments

Presenter: **Jiayi Wang**, The University of Texas at Dallas, United States

Co-authors: Zhengling Qi, Chengchun Shi

As AI becomes more prevalent throughout society, effective methods of integrating humans and AI systems that leverage their respective strengths and mitigate risk have become an important priority. The paradigm of super reinforcement learning is introduced that takes advantage of human-AI interaction for data-driven sequential decision making. This approach utilizes the observed action, either from AI or humans, as input for achieving a stronger oracle in policy learning for the decision maker (humans or AI). In the decision process with unmeasured confounding, the actions taken by past agents can offer valuable insights into undisclosed information. By including this information for the policy search in a novel and legitimate manner, the proposed super reinforcement learning will yield a super-policy that is guaranteed to outperform both the standard optimal policy and the behavior one (e.g., past agents' actions). This stronger oracle is called a blessing from human-AI interaction. Furthermore, to address the issue of unmeasured confounding in finding super-policies using the batch data, a number of nonparametric and causal identifications are established. Building upon these novel identification results, several super-policy learning algorithms are developed, and their theoretical properties are systematically studied, such as finite-sample regret guarantee.

C1454: Pairwise personalized learning in recommendation system

Presenter: **Haowen Zhou**, University of Virginia, United States

Co-authors: Xiwei Tang, James Lee

Recommendation systems are essential across domains such as content platforms and e-commerce, yet standard approaches based solely on explicit ratings often face challenges due to data sparsity and heterogeneity. The aim is to introduce a hybrid pairwise personalized learning (HPPL) model that extends the conventional Bayesian personalized ranking method to handle a mixture of explicit and implicit feedback in recommendation tasks. Unlike traditional models that treat observed data as weakly positive and missing data as weakly negative signals, HPPL leverages explicit ratings in a weighted pairwise loss function with latent factors, achieving robust performance in ranking-based evaluations while maintaining computational efficiency. The model bridges the gap between explicit and implicit feedback systems, providing theoretical guarantees for scalable and adaptive stochastic gradient algorithms and offering a practical solution for real-world recommendation systems.

CO193 Room BCB G08 BAYESIAN TIME-SERIES MODELLING

Chair: Roberto Casarin

C0684: A dynamic stochastic block model for multidimensional networks

Presenter: **Ovielt Antonio Baltodano Lopez**, Ca' Foscari University, Italy

Co-authors: Roberto Casarin

The availability of relational data can offer new insights into the functioning of the economy. Nevertheless, modeling the dynamics in network data with multiple types of relationships is still a challenging issue. Stochastic block models provide a parsimonious and flexible approach to network analysis. A new stochastic block model is proposed for multidimensional networks, where layer-specific hidden Markov-chain processes drive the changes in community formation. The changes in the block membership of a node in a given layer may be influenced by its own past membership in other layers. This allows for clustering overlap, clustering decoupling, or more complex relationships between layers, including settings of unidirectional or bidirectional, non-linear Granger block causality. The overparameterization issue of a saturated specification is coped with by assuming a Multi-Laplacian prior distribution within a Bayesian framework. Through simulations, it is shown that standard linear models and the pairwise approach are unable to detect block causality in most scenarios. In contrast, the model can recover the true Granger causality structure. As an application to international trade, it is shown that the model offers a unified framework, including community detection and the Gravity equation modeling. New evidence of block Granger causality of trade agreements and trade flows is found.

C0685: Bayesian outlier detection for matrix-variate models

Presenter: **Antonio Peruzzi**, Ca' Foscari University of Venice, Italy

Co-authors: Roberto Casarin, Monica Billio, Fausto Corradin

Bayes factor (BF) is one of the tools used in Bayesian analysis for model selection. The predictive BF finds application in detecting outliers, which are relevant sources of estimation and forecast errors. An efficient framework for outlier detection is provided and purposely designed for large multidimensional datasets. Online detection and analytical tractability guarantee the procedure's efficiency. The proposed sequential Bayesian monitoring extends the univariate setup to a matrix-variate one. Prior perturbation based on power discounting is applied to obtain tractable predictive BFs. This way, computationally intensive procedures used in Bayesian analysis are not required. The conditions leading to inconclusive responses in outlier identification are derived, and some robust approaches are proposed that exploit the predictive BF's variability to improve the standard discounting method. The effectiveness of the procedure is studied using simulated data. An illustration is provided through applications to relevant benchmark datasets from macroeconomics and finance.

C0686: Bayesian compressed tensor regression

Presenter: **Qing Wang**, Ca Foscari University, Italy

Co-authors: Roberto Casarin, Radu Craiu

A new dimensionality reduction technique is proposed to address the common problem of high dimensionality in tensor regressions. A generalized tensor random projection method is introduced that embeds high-dimensional tensor-valued covariates into low-dimensional subspaces with minimal loss of information about the responses. The method is flexible, allowing for tensor-wise, mode-wise, or combined random projections as special cases. A Bayesian inference framework is provided featuring the use of a hierarchical prior distribution and a low-rank representation of the parameter. Strong theoretical support is provided for the concentration properties of the random projection and posterior consistency of the Bayesian inference. An efficient Gibbs sampler is developed to perform inference on the compressed data. To mitigate the sensitivity introduced by random projections, Bayesian model averaging is employed, with normalizing constants estimated using reversed logistic regression. An extensive simulation study is conducted to examine the effects of different tuning parameters. A real data application demonstrates that compressed Bayesian tensor regression achieves better out-of-sample prediction while significantly reducing computational cost compared to standard Bayesian tensor regression.

C0757: An observation-driven generalized Poisson model

Presenter: **Dario Palumbo**, University Ca Foscari of Venice, Italy

Co-authors: Roberto Casarin, Giulia Carallo

Based on the generalized Poisson (GP) conditional distribution, a new general class of observation-driven models for count data is presented, and their theoretical properties are derived. The GP is a flexible distribution that allows for both under- and over-dispersion. As a special member of this class, a score-driven model version is introduced. It is shown that this specification is robust to the presence of outliers and can be extended to allow for time-varying over-dispersion. For the estimation of the model, a Bayesian inference framework and an efficient posterior approximation procedure based on Markov Chain Monte Carlo are provided. Posterior contraction rates are also established with increasing sample size in terms of the average Hellinger metric. The applications on environmental variables show that the proposed model is well-suited for capturing the over-dispersion feature of the data.

C0854: Heterogeneous g-priors for networks with dyadic covariance structures

Presenter: **Rigers Behluli**, University of Warwick - Ca'Foscari University, Italy

Co-authors: Roberto Casarin, Mark Steel

The increasing availability of multivariate data calls for the use of matrix variate models for identifying hidden patterns within the data and for predicting the variables of interest. The aim is to propose a novel regression model for sequences of matrix data with dyadic covariance structure and develop a Bayesian procedure for model selection. It is shown that it is possible to characterize the dyadic variance-covariance matrix analytically and perform inference and model selection using mixtures of g-priors adapted to accommodate the heterogeneity induced by dyadic covariance. Some theoretical results are presented, and the algorithms used to sample low- and high-dimensional model spaces. Simulation results and a real data comparison to an established method for dyadic covariance estimation are presented.

CO075 Room BCB G09 TOPICS IN FINANCIAL ECONOMETRICS

Chair: Leopold Soegner

C0274: Forecasting inflation with the hedged random forest

Presenter: **Michael Wolf**, University of Zurich, Switzerland

Co-authors: Elliot Beck

Accurately forecasting inflation is critical for economic policy, financial markets, and broader societal stability. In recent years, machine learning methods have shown great potential for improving the accuracy of inflation forecasts; specifically, random forests stand out as a particularly effective approach that consistently outperforms traditional benchmark models in empirical studies. Building on this foundation, the hedged random forest (HRF) framework of a prior study is adapted for the task of forecasting inflation. Unlike the standard random forest, the HRF employs non-equal (and even negative) weights of the individual trees, which are designed to improve forecasting accuracy. Estimators of the HRF's two inputs, the mean and the covariance matrix of the errors corresponding to the individual trees, that are customized for the task at hand, are developed. An extensive empirical analysis demonstrates that the proposed approach consistently outperforms the standard random forest.

C0413: Machine learning forecasting of industrial production in Slovakia

Presenter: **Adam Csapai**, The Institute of Economic Research of Slovak Academy of Sciences and University of Economics in Bratislava, Slovakia

An extensive evaluation applies state-of-the-art machine learning models to the task of forecasting industrial production in Slovakia, demonstrating their superior ability to capture the evolution of a small, open, industrialized economy operating within a monetary union. The purpose is to represent the first systematic assessment of machine learning forecasting performance in such an environment characterized by short time series and two crisis episodes, in contrast to prior work on larger economies with decades of data. A key finding is that regularization serves as a more effective dimension-reduction technique than principal component analysis or factor models, preserving essential information and yielding more accurate predictions. Forecast combination methods that weight models according to both error magnitude and directional accuracy are also applied, and the directional accuracy of machine learning methods is explored. The robustness of these techniques is tested across pre- and post-COVID-19 periods to evaluate model performance under heightened volatility and uncertainty. Finally, an analysis of soft indicators as standalone inputs highlights their limited predictive value, and earlier critiques are also addressed by underscoring the practical advantages of machine learning approaches in macroeconomic forecasting.

C1067: Online breakpoint-detection in cointegrating relationships

Presenter: **Leopold Soegner**, Institute for Advanced Studies, Austria

Co-authors: Martin Wagner

A closed-end consistent monitoring procedure is developed with the goal of detecting structural changes in cointegrating relationships. The vector error correction model is considered, and different specifications of the deterministic terms are allowed for. βy_t is considered, where β is a matrix containing the cointegrating vectors and (y_t) is a process integrated of order one, and propose a monitoring test statistic to investigate the stability of these cointegrating relationships. The asymptotic distribution of the test statistic is obtained under the null hypothesis of no structural breaks. A calibration period is used for parameter estimation, after which online break-point detection is performed. The procedure stops at the first time point the test statistic exceeds the corresponding critical value. It is shown that the monitoring procedure is consistent. A simulation study is provided to investigate the finite sample properties of our monitoring procedure. The procedure is applied to investigate the triangular exchange rate parity.

C1092: Fully modified estimation of a quantile cointegration model with a spatial lag

Presenter: **Christian Haefke**, New York University, United Arab Emirates

Co-authors: Leopold Soegner

The purpose is to investigate a quantile cointegration regression problem including a spatial lag. The asymptotic distribution of the quasi-maximum likelihood estimator is obtained, and it is shown that the second-order bias depends on the order of the deterministic terms. Prior studies are followed to construct a fully modified estimator, which removes this second-order bias and a Wald-type test. The finite sample properties of the fully modified estimator are analyzed in a simulation study. The relevance of the estimator is assessed by comparing empirical COVAR estimates of a prior study.

C1099: A stock-flow consistent analysis of credit cycles and the productivity puzzle

Presenter: **Saite Lu**, Emmanuel College, University of Cambridge, United Kingdom

The purpose is to investigate the macroeconomic consequences of credit cycles, focusing on the role of mortgage-driven household consumption in shaping growth and productivity dynamics. Employing a stock-flow consistent (SFC) model grounded in Nicholas Kaldors growth and productivity hypothesis, it is examined how household credit expansions, particularly through mortgage borrowing, affect aggregate demand and long-run economic performance. Kaldor's hypothesis posits that productivity growth is not exogenously given but endogenously driven by output growth, especially in the presence of increasing returns to scale. Accordingly, rising household credit initially stimulates output and productivity through enhanced demand, but over time, persistent reliance on debt-financed consumption leads to rising financial fragility, weaker investment, and slow growth. These dynamics offer a structural explanation for the productivity slowdown observed in many advanced economies following the global financial crisis (GFC). The model captures key interactions between households, firms, banks, and the government. Simulations calibrated to stylized macroeconomic trends reveal how credit cycles amplify macroeconomic instability and erode productivity gains. The contribution to the literature is on the post-GFC productivity puzzle by linking financial structures and household behavior to demand-led growth theory, and by underscoring the regulatory role of fiscal and macroprudential policy.

CO139 Room Virtual R01 STATISTICAL ADVANCES IN MACHINE LEARNING FOR COMPLEX BIOMEDICAL DATA Chair: Li-Xuan Qin**C0767: Contextual evaluation of data harmonization for microRNA sequencing***Presenter:* **Li-Xuan Qin**, Memorial Sloan Kettering Cancer Center, United States

The reproducibility of microRNA sequencing data analysis hinges on effectively mitigating data artifacts that arise from variable experimental handling through data harmonization. While numerous harmonization methods encompassing normalization and batch-effect correction have been developed to address these artifacts, statistical investigations into their impact on downstream analyses primarily focused on differential expression analysis. To enable contextual evaluation for data harmonization, robust benchmark datasets are developed, thorough evaluation pipelines, and accompanying software tools, with a particular focus on microRNAs. Findings are presented from a simulation study evaluating the performance of various data harmonization approaches in the contexts of sample clustering and sample classification, each assessed using multiple analytical methods. The best-performing combinations of harmonization and downstream analysis methods were then applied to reanalyze publicly available real-world data.

C0877: Variable selection in compositional data analysis*Presenter:* **Jing Ma**, Fred Hutchinson Cancer Center, United States*Co-authors:* Kristyn Pantoja, David Jones

Compositional data, where only relative abundances are available, are common in microbiome and other high-throughput sequencing studies. Log ratios between groups of variables serve as key biomarkers in these settings. However, selecting predictive log ratios is a combinatorial challenge, and existing greedy search-based methods are computationally expensive, limiting their applicability to high-dimensional data. The supervised log ratio (SLR) method is introduced, a novel and efficient approach for selecting predictive log ratios in high-dimensional settings. SLR first screens active variables using univariate regression on log ratio transformed data and then applies principal balance analysis to define balance biomarkers. The approach leverages both the relationship between the response and predictors and the correlations among the predictors to improve accuracy in variable selection and prediction. Through simulations and two case studies, one on inflammatory bowel disease (IBD) and another on colorectal cancer (CRC), it is demonstrated that SLR outperforms existing methods, particularly in high-dimensional settings.

C0781: Interactively resolving distortion in nonlinear dimensionality reduction of biomedical data*Presenter:* **Kris Sankaran**, University of Wisconsin, United States

Nonlinear dimensionality reduction is a key step in many biomedical analysis workflows. For example, when working with text embeddings from pretrained protein language models or when exploring single-cell gene expression measurements, researchers routinely apply UMAP to organize the high-dimensional source data into a more manageable low-dimensional representation. Such nonlinear dimensionality can be powerful, but it inevitably introduces distortion. A growing body of work has demonstrated that this distortion can have serious consequences for downstream interpretation, for example, suggesting clusters that do not exist in the original data. Motivated by these developments, a visual interface is designed that helps to identify where these distortions are most severe and supports interaction to locally resolve them. Though the design and interaction are relatively straightforward, it is found through case studies from single-cell genomics and microbiome data analysis that they can enable more accurate interpretations than more traditional visualization methods, which do not show distortion. It helps researchers who apply nonlinear dimensionality reduction methods address concerns they may have about the reliability of their embeddings and proceed with confidence in their data analysis.

C1028: Error-controlled non-additive interaction discovery in machine learning models*Presenter:* **Yang Lu**, University of Wisconsin-Madison, United States

Machine learning (ML) models are powerful tools for detecting complex patterns, yet their "black box" nature limits their interpretability, hindering their use in critical domains like healthcare and finance. Interpretable ML methods aim to explain how features influence model predictions but often focus on univariate feature importance, overlooking complex feature interactions. While recent efforts extend interpretability to feature interactions, existing approaches struggle with robustness and error control, especially under data perturbations. Diamond, a method for trustworthy feature interaction discovery, is introduced. Diamond uniquely integrates the model-X knockoffs framework to control the false discovery rate (FDR), ensuring a low proportion of falsely detected interactions. Diamond includes a non-additivity distillation procedure that refines existing interaction importance measures to isolate non-additive interaction effects while preserving FDR control. This approach addresses the limitations of off-the-shelf interaction measures, which, when used naively, can lead to inaccurate discoveries. Diamond's applicability spans a broad class of ML models. Empirical evaluations on both simulated and real datasets across various biomedical studies demonstrate its utility in enabling reliable data-driven scientific discoveries. Diamond represents a significant step forward in leveraging ML for scientific innovation and hypothesis generation.

C0961: Integrated analysis of imaging and RNA-seq to characterize smoking-related lung disease phenotypes*Presenter:* **Fenghai Duan**, Brown University School of Public Health, United States

The goal is to identify biological distinctions across smoking-related phenotypes by integrating clinical, imaging, and bronchial epithelial RNA-seq data. Using k-means clustering, participants are grouped based on CT imaging features, and their clinical phenotypes are analyzed. Bronchial epithelial RNA-seq data are further examined through differential gene expression, gene set enrichment, and variation analyses. Three distinct clusters are identified: Preserved, interstitial predominant, and emphysema predominant. Compared to the preserved cluster, the interstitial and emphysema clusters showed poorer lung function, lower exercise capacity, and worse quality of life. They also experienced faster declines in function, greater emphysema progression, more respiratory events, and higher mortality. The emphysema cluster had the most severe outcomes, followed by the interstitial cluster. Transcriptomic analysis indicated that severe disease stages were linked to heightened inflammatory responses, especially through the TNF-alpha pathway, while milder stages showed upregulation of T-cell-related genes. Using quantitative CT imaging, we identified three subgroups among individuals with a history of heavy smoking. Differences in airway gene expression suggest a connection between clinical severity and inflammation, possibly mediated by the TNF-alpha pathway.

CO278 Room BCB 206 CFE SESSION: A TRIBUTE TO H. PESARAN**Chair: Ilias Chronopoulos****C0659: Panel regressions with latent factors and two-way heterogeneous treatments: A unified perspective***Presenter:* **Martin Weidner**, University of Oxford, United Kingdom*Co-authors:* Arturas Juodis

Panel regressions are revisited with unobserved heterogeneity through the lens of variance-weighted average treatment effects. Building on established results for cross-sectional OLS and one-way fixed effects panels, it is shown that two-way panel estimators with latent factors, such as the interactive fixed effects estimator, also converge to interpretable estimands under minimal assumptions. Specifically, even when the factor structure is misspecified, the estimator targets a variance-weighted average of unit-time-specific treatment effects, generalizing insights from prior studies. This unified perspective bridges the treatment effects literature and the factor regression literature, highlighting which estimators admit meaningful interpretation in heterogeneous panels without relying on parametric restrictions. Implications are also outlined for empirical practice and suggest directions for robust inference under latent interactive structures.

C0691: Estimation and inference in large dimensional threshold factor models with weaker loadings*Presenter:* **Daniele Massacci**, Kings College London, United Kingdom

Estimation and inference are studied in large-dimensional threshold factor models in which (some of) the eigenvalues of the covariance matrix

of the data diverge at a rate that is slower than the cross-sectional dimension N . The convergence rate of the concentrated least squares estimator is derived for the threshold parameter, and the asymptotic distribution of the principal components estimator is obtained for factors and loadings within each regime. Finally, the relevance of these theoretical findings is illustrated for conditional asset pricing.

C1127: Incorporating micro data into macro models using pseudo VARs

Presenter: **James Mitchell**, Federal Reserve Bank of Cleveland, United States

Co-authors: Gary Koop, Stuart McIntyre, Ping Wu

The purpose is to develop a method to incorporate microdata, available as repeated cross-sections, into macro VAR models to understand the distributional effects of macroeconomic shocks. The method extends existing functional VAR models by “looking within” the micro distribution to identify the degree to which specific types of micro agents are affected by the shock. It does so by creating a pseudo-panel from the repeated cross-section and adding these pseudo-individuals into the macro VAR. Jointly modeling the micro and macro data leads to a large (pseudo) VAR. Bayesian methods are used to ensure shrinkage and parsimony. The application revisits a prior study and compares their functional VAR-based distributional impulse response functions with the proposed pseudo VAR-based ones to identify which types of individuals’ earnings are most affected by business cycle shocks.

C1308: Model selection in large dimensional linear regression using sequential multiple testing

Presenter: **Vasilis Sarafidis**, Brunel University London, United Kingdom

Co-authors: George Kapetanios, Alexia Ventouri

High-dimensional regression specification and analysis is a complex and active area of research in statistics and econometrics. A large number of approaches have been proposed, but determining their relative merits remains a challenging task. The aim is to propose a new hybrid approach. It combines elements from two existing methods. The first is the greedy methodology developed in a prior study, where a powerful multiple testing step introduces parsimony by ensuring that completely irrelevant variables are not selected with high probability. The second is stage-wise regression, where relevant variables are selected in steps, rather than jointly as in the prior study. However, in that literature, the stopping rules, which typically rely on model information criteria, are not sufficiently parsimonious. Some theoretical properties of the new method are derived and shown, through simulations, to perform well. An illustration, using corporate emissions data, provides an empirical perspective.

CO232 Room BCB 207 MACHINE LEARNING METHODS FOR NOWCASTING

Chair: Eric Ghysels

C0379: An observation-driven mixed-frequency VAR model

Presenter: **Heiner Mikosch**, ETH Zurich, Switzerland

Co-authors: Maurizio Daniele, Stefan Neuwirth

The aim is to introduce an observation-driven mixed-frequency vector autoregression (MFVAR) model. A general MFVAR framework is developed based on vector stacking, and it is shown how the model can be reformulated into a closed-form representation that permits analytical estimation. Additionally, a Bayesian normal prior is derived to enable shrinkage estimation of the MFVAR. Monte Carlo simulations and empirical applications demonstrate that the estimators are computationally efficient, even in high-dimensional settings. The observation-driven MFVAR achieves comparable or superior out-of-sample forecasting accuracy relative to state-space MFVARs, while requiring only a fraction of the computation time.

C0642: Nowcasting and aggregation: Why small Euro area countries matter

Presenter: **Luca Barbaglia**, European Commission Joint Research Centre, Italy

Co-authors: Andrii Babii, Eric Ghysels, Jonas Striaukas

The purpose is to study the nowcasting of Euro area Gross Domestic Product (GDP) growth using mixed data sampling machine learning panel data regressions with both standard macro releases and daily news data. Using a panel of 19 Euro area countries, it is investigated whether directly nowcasting the Euro area aggregate is better than weighted individual country nowcasts. Results highlight the importance of the information from small- and medium-sized countries, particularly when including the COVID-19 pandemic period. The analysis is supplemented by studying the so-called Big Four (France, Germany, Italy, and Spain) and the value added of news data when official statistics are lagging.

C0790: Nowcasting with functional approaches and mixed-frequency data

Presenter: **Anna Simoni**, GENES - CREST, France

Co-authors: Catherine Doz, Laurent Ferrara

Prediction and causal models often involve economic time series with different sampling frequencies. The mix-frequency problem has received a lot of attention in the past nowcasting/forecasting literature, which has proposed solutions that work especially well for prediction when the frequency gap is small, like the gap between monthly and quarterly data. An alternative approach is proposed to deal with mix-frequency, which is particularly designed for situations where the gap in sampling frequencies is large, like daily and quarterly. Such a large gap can be easily found when one mixes series from standard and non-standard sources, like the internet or newspapers. The approach focuses on both prediction and causality, and exhibits excellent performance. By treating the high-frequency variable as a realization of a stochastic process in continuous time, the estimation problem is cast in the class of ill-posed inverse problems, and different regularized estimators are proposed. Importantly, for each high-frequency covariate, a measure of its influence on the target variable is recovered as a function of the time-gap between the forecasting horizon and the date of the information in the past. Because the analysis is conditional on several covariates, the fact that the influence of high-frequency series rapidly decreases, as soon as information from standard data arises, is well-captured.

C0221: LLMs vs econometric models for nowcasting GDP growth: A practitioner’s view

Presenter: **Marie Bessec**, University Paris Dauphine, France

Co-authors: Julien Andre, Zachary Goulby

The performance of large language models (LLMs) is evaluated in nowcasting French GDP growth, comparing them with the econometric models currently used by the Banque de France. Using only prompt-based queries without external data or fine-tuning, it is assessed whether general-purpose LLMs such as ChatGPT, Gemini, and Claude can serve as effective forecasting tools. While econometric models consistently outperform LLMs during normal periods, the latter show a notable advantage in capturing exceptional events such as the COVID-19 pandemic. The sensitivity of LLM forecasts is also examined to prompt language, design, and model version, and a confidence index and recession probability derived from LLM responses are introduced. There is no strong evidence of information leakage, and robustness checks confirm the findings across various model versions and temperature settings. A fair in-sample comparison reinforces the relative strength of econometric models in normal conditions. Overall, the results suggest that while standard LLMs are not yet ready to replace traditional models in routine forecasting, they can provide complementary insights, particularly in periods of structural change or heightened uncertainty. These results apply to the most popular LLMs without external data and fine-tuning.

CO241 Room BCB 209 NONCAUSAL TIME SERIES

Chair: Sean Telg

C0434: Tail-aware density forecasting of locally explosive time series: A neural network approach

Presenter: **Julien Peignon**, Paris Dauphine University, France

Co-authors: Arthur Thomas, Elena Dumitrescu

Mixed-causal ARMA processes are known to capture the dynamics of locally explosive behavior, such as bubble assets in finance. However, the limited knowledge of the predictive density of mixed-causal processes, especially during explosive bubble events, complicates their forecast and thus limits their use in practical applications. Given the lack of closed-form formulae for the conditional prediction density (except in special cases), simulation-based and sample-based methods have been proposed in the literature. However, these methods can be computationally expensive and do not accurately capture the dynamics during explosive episodes. Mixture density networks (MDNs) are introduced for forecasting time series that exhibit locally explosive behavior. By incorporating Tukey g-and-h transformations as mixture components, the approach offers enhanced flexibility in capturing the skewed, heavy-tailed, and potentially multimodal nature of predictive densities associated with bubble dynamics. Furthermore, the weighted likelihood emphasizes tail observations and crash events, enabling accurate density estimation in the extreme regions most relevant for risk management. Finally, once trained, the MDN produces near-instantaneous density forecasts. Through extensive Monte Carlo simulations and empirical applications, it is shown that the proposed MDN-based framework delivers superior forecasting performance relative to existing approaches.

C0473: Deconvolution and filtering of non-causal alpha-stable processes

Presenter: **Ludivine Vaudree**, University of Orleans, France

Co-authors: Gilles De Truchis, Arthur Thomas

The purpose is to develop a comprehensive theoretical framework for the deconvolution and filtering of time series processes composed of aggregated non-causal alpha-stable AR(1) components. Rigorous identification conditions are established based on the characteristic function, and efficient sequential estimation methods are proposed. The approach allows for the modeling of multiple local explosive behaviors occurring at different rates, such as those observed in financial time series exhibiting speculative bubbles. Upon establishing parameter identification, a particle filtering method is introduced based on Markov chain Monte Carlo techniques, specifically designed for the recovery of latent components in alpha-stable mixture models. This filtering methodology effectively addresses the challenges inherent to alpha-stable distributions, including heavy tails and the lack of closed-form densities. The framework is developed for both continuous and discrete domains. This combination of parameter estimation and filtering techniques enables the analysis of financial time series exhibiting locally explosive patterns while preserving stationarity.

C0889: Weighted maximum likelihood for misspecified mixed causal non-causal autoregressive models

Presenter: **Gabriele Mingoli**, Vrije Universiteit, Netherlands

Co-authors: Francisco Blasques, Siem Jan Koopman

The aim is to introduce a novel weighted maximum likelihood estimation (WMLE) method aimed at improving the forecast accuracy of misspecified models. The approach is specifically designed for mixed causal non-causal autoregressive (MAR) models, which allow a stochastic process to depend on its future values through a lead polynomial, capturing locally explosive dynamics. MAR models are commonly used to describe and predict speculative bubbles in financial time series. However, they impose a rigid structure on all bubbles in a given dataset, implying that all bubbles exhibit the same dynamics, which may not always align with empirical observations. In reality, explosive episodes in time series often vary in size, duration, and growth rate, making the standard MAR specification restrictive. It is demonstrated that by applying different weights to different parts of the sample when estimating a MAR model, it is possible to construct an estimator that enhances forecasting performance compared to traditional maximum likelihood estimation (MLE). A simulation study confirms that the proposed weighting approach improves predictive accuracy under misspecification. Finally, an out-of-sample forecasting exercise on monthly crude oil prices shows that WMLE outperforms forecasts obtained using standard MLE.

C0900: A generalized methods of moments approach for a noncausal dynamic panel model

Presenter: **Kevin Cecere Palazzo**, Vrije Universiteit, Netherlands

Co-authors: Sean Telg, Siem Jan Koopman, Francisco Blasques

A generalized method of moments approach is proposed for estimation and identification of a noncausal dynamic panel model. The efficiency of the proposed estimator is investigated through the inclusion of additional instruments compared to standard references in the literature on the causal dynamic panel model. A panel framework permits more powerful tools for estimation and identification. A simulation exercise investigates the behavior of the proposed estimator for different assumptions on the data-generating process.

C0936: Introducing fractional integration in mixed causal-noncausal models

Presenter: **Sean Telg**, Vrije Universiteit Amsterdam, Netherlands

Co-authors: Sebastien Fries, Jorik van der Oord, Jean-Michel Zakoian

The notion of fractional integration is introduced into the mixed causal-noncausal autoregressive (MAR) model. It is shown that this new model, called FIMAR, is able to generate baseline paths that exhibit combinations of exponential and hyperbolic growth as often observed in speculative bubbles. For a general class of error distributions, stationarity conditions and the existence of a purely causal second-order equivalent (SOE) form of the FIMAR model are derived. Using this SOE representation in combination with a Whittle-type estimator, it is demonstrated how to perform model selection for FIMAR models.

CO171 Room BCB 210 ADVANCES IN TIME SERIES FOR ECONOMICS AND FINANCE

Chair: Luca Scaffidi Domianello

C0276: Trend in Markov-switching VAR models

Presenter: **Maddalena Cavicchioli**, University of Modena and Reggio Emilia, Italy

The purpose is to consider Markov-switching vector autoregressive (MS VAR) processes that incorporate both a stationary component and a deterministic time trend. These models are particularly relevant for capturing structural changes and asymmetries in macroeconomic time series. To estimate the model parameters, the ordinary least squares (OLS) method is adopted within a modified expectation-maximization (EM) algorithm framework. This approach offers computational simplicity while accommodating regime changes governed by an unobserved Markov process. The consistency is established, and the asymptotic distribution of the resulting OLS estimators is derived. To the best of knowledge, the characterization of the asymptotic variance-covariance matrix for OLS estimators in this context fills a gap in the existing econometric literature. An empirical application focused on housing market asymmetries demonstrates the practical relevance and effectiveness of the proposed methodology.

C0418: Combining GARCH-MIDAS forecasts of US state-level volatility: The role of local and global EPU indices

Presenter: **Vincenzo Candila**, University of Salerno, Italy

Co-authors: Oguzhan Cepni, Giampiero Gallo, Rangan Gupta

The role of both local (state-specific) and global economic policy uncertainty (EPU) is investigated in forecasting the volatility of U.S. state-level equity returns. A GARCH-MIDAS framework is adopted, which allows the inclusion of multiple EPU indices as low-frequency drivers of daily return volatility. To tackle the issue of identifying the most relevant predictors, an elastic net (EN) shrinkage technique is implemented that combines forecasts across different model specifications. Findings show that the combined model, which leverages information from both local and global EPU indices, consistently outperforms individual specifications, from both one-step- and multi-step-ahead perspectives. These results underscore the importance of accounting for both regional and global policy uncertainty when modeling volatility dynamics at the state level.

C0421: A structural matrix autoregressive model for volatility, volume, jumps and returns

Presenter: **Andrea Bucci**, University of Macerata, Italy

Co-authors: Giulio Palomba, Eduardo Rossi

The aim is to investigate commonalities between daily traded volumes, returns, volatility, and jumps. The approach relies on a structural matrix-variate model to potentially account for cross-variables and cross-asset spillovers simultaneously. A full identification is provided for the cross-variables structural relationships based on the financial theory, and the model is applied to 50 assets from the S&P 500 index. The results indicate that jumps positively affect the other variables and that volatility severely affects trading volumes.

C0471: Matrix-variate hidden Markov models for robust clustering and anomaly detection

Presenter: **Salvatore Daniele Tomarchio**, University of Catania, Italy

The matrix-variate hidden Markov model (HMM) framework is extended by introducing two novel model families that employ the matrix-variate t distribution and a contaminated normal distribution. These extensions enhance the modeling of heavy tails, improve clustering performance, and support the detection of atypical matrices. Parameter estimation relies on two tailored expectation-conditional maximization (ECM) algorithms, both of which are implemented in the MatrixHMM R package. Simulation studies demonstrate the models' accuracy, robustness, and effectiveness in outlier identification. Finally, an application to real data illustrates how the models can be used to explore labor market trends across Italian provinces.

C1217: Common features in volatilities: A new multiplicative error model

Presenter: **Luca Scaffidi Domianello**, University of Catania, Italy

Co-authors: Edoardo Otranto

Multivariate volatility models could consider the influence of each volatility series on the others (spillover effects). Furthermore, integrating financial markets provides similar dynamics (co-movements). The aim is to propose a new model for volatility vectors, belonging to the family of multiplicative error models (MEMs), which incorporates spillover and co-movement effects captured through a separate unobservable component. Moreover, to reduce the number of coefficients for high-dimensional datasets, a simple model-based clustering procedure is proposed. The model is applied to a set of 29 assets included in the Dow Jones Industrial Index, providing insights into the interpretation of spillover effects and co-movement. The adopted parametrization shows a satisfactory performance when compared to other vector MEMs.

CO018 Room BCB 211 ECOSTA JOURNAL SESSION

Chair: Simona Sanfelici

C1455: Tractable unified skew- t distribution and copula for heterogeneous asymmetries

Presenter: **Michael Smith**, University of Melbourne, Australia

Co-authors: Lin Deng, Worapree Ole Maneeoonthorn

Multivariate distributions that allow for asymmetry and heavy tails are important building blocks in many statistical models. The unified skew- t (UST) is a promising choice because it is scalable and allows for a high level of flexibility in the asymmetry of distribution. However, it suffers from parameter identification and computational hurdles that have to date inhibited its use for modeling data. The aim is to propose a new tractable variant of the unified skew- t (TrUST) distribution that addresses both challenges. Moreover, the copula of this distribution is shown to also be tractable, while allowing for greater heterogeneity in asymmetric dependence over variable pairs than the popular skew- t copula. It is shown how Bayesian posterior inference for both the distribution and its copula can be computed using an extended likelihood derived from a generative representation of the distribution. The efficacy of this Bayesian method, and the enhanced flexibility of both the TrUST distribution and its implicit copula, is first demonstrated using simulated data. Applications of the TrUST distribution to highly skewed Australian electricity prices, and the TrUST copula to intraday U.S. equity returns, demonstrate how the proposed distribution and its copula can provide substantial increases in accuracy in practice.

C1305: A two-sample smooth test for dependent data

Presenter: **Eric Beutner**, Vrije Universiteit Amsterdam, Netherlands

The purpose is to consider a two-sample smooth test for testing the equality of multivariate distributions. Dependency between the two samples is allowed for. Consistency of the two-sample smooth test for dependent samples is shown. Moreover, various bootstrap schemes are considered for this two-sample smooth test, and bootstrap validity is established.

C1145: Junction tree structured Markov random fields with Bernoulli marginals

Presenter: **Etienne Marceau**, Laval University, Canada

Co-authors: Helene Cossette, Anthony Gelin

The aim is to investigate the family of junction tree structured Markov random fields with Bernoulli marginal. This allows for a generalization of the tree-based Ising model with improved flexibility in terms of dependence modeling, while retaining the benefits of mean parametrization. An analytic expression is provided for the joint probability mass function of any Bernoulli random vector using moments as parameters, similar to other results found in the literature. This allows for an analytic form for the probability mass function of the studied model. An analytic expression is derived for the joint probability generating function of Bernoulli random vectors encrypted on junction tree structured Markov random fields, along with some applications and algorithms for its computation. The properties of junction trees also allow for a clique-based representation, which yields an efficient sampling algorithm for the model. The estimation of the model is finally assessed from data for a set maximal clique size with maximum likelihood estimation and an analysis on the complexity of the model. The necessary algorithms are provided to estimate models and a numerical example using precipitation data in a high dependence context.

C1260: Accelerating estimation for GARCH-type models by weighted data subsampling

Presenter: **Matias Quiroz**, University of Technology Sydney, Australia

Co-authors: Zixuan Wang, Aishwarya Bhaskaran, Thomas Goodwin

Generalised autoregressive conditional heteroskedasticity (GARCH) models are ubiquitous in financial econometrics to capture and forecast time-varying volatility in asset returns. Thus, they are essential tools in the econometrician's toolbox for risk management, derivative pricing, and portfolio optimisation, and beyond. Likelihood-based estimation of GARCH models is computationally expensive due to the recursive structure of the conditional variance, limiting their applicability in large datasets and scenarios requiring quick decision-making. To speed up the estimation, a subsampling-based unbiased log-likelihood estimator is proposed that includes early observations with a higher probability, thereby reducing the length of the recursive conditional variance loop when evaluating the likelihood. A subsampling-based unbiased gradient of the log-likelihood estimator is also proposed. The proposed estimators may be used in any subsampling-based inferential approach, such as stochastic optimization to find the maximum likelihood estimate in the classical paradigm, or subsampling Markov chain Monte Carlo (MCMC) algorithms to sample the posterior distribution in the Bayesian paradigm.

C0230: Fourier-Malliavin volatility estimation: Recent advances and applications

Presenter: **Simona Sanfelici**, University of Parma, Italy

The Fourier volatility estimation method has since stimulated a growing body of scientific literature. Numerous papers have explored both its theoretical foundations and practical applications, published in peer-reviewed international journals across diverse fields - including mathematical finance, high-frequency econometrics, econophysics, and even extending to the natural and medical sciences. The purpose is to present a selection of recent theoretical advancements and updated financial applications of the Fourier estimator. Owing to its ability to reconstruct volatility as a stochastic function of time in both univariate and multivariate settings, the Fourier-Malliavin methodology offers deep insights into a range

of volatility-related financial quantities. These include volatility of volatility, leverage effects, and other second-order effects such as the price-volatility feedback rate. This line of research has continued to evolve in recent years. Recent contributions are highlighted, particularly relevant to early warning systems for detecting financial instability, and to factor identification in stochastic volatility models.

CO050 Room BCB 212 ADVANCES IN RISK MEASUREMENT AND IN FORECASTING
Chair: Roxana Halbleib
C0303: Factor inference under common components in volatility

Presenter: **Julia Koh**, Tilburg University, Netherlands

Co-authors: Benoit Perron, Silvia Goncalves

Considering the common components of volatility is crucial in finance. It is found that the standard factor literature does not account for these common components in volatility. Alternative assumptions are introduced that can accommodate the volatility dependence and develop inference methods for factor models. Three key applications are revisited: Factor estimation and distribution theory, factor-augmented regression models, and a test for the number of factors in group factor models. Under new assumptions, it is shown that the results for factor estimation and factor-augmented regression models are maintained. It is found that the variance of the test statistic for the number of factors in the group model is affected, and a new estimator for the variance is proposed.

C0368: Combining forecasts based on the evidence against equal weights

Presenter: **Lukas Bauer**, University of Freiburg, Statistics and Econometrics, Germany

A novel class of performance-based forecast combination schemes is proposed. The schemes determine the combining weights using the standardized loss difference of each model relative to the average model, and thus account for the statistical magnitude of the differences in predictive ability. The risk of two such combining schemes is characterized relative to the best individual model. These schemes build an intuitive bridge between model selection and forecast combination while showing robust performance in finite samples. An empirical application is performed, combining forecasts of financial risk, i.e., realized volatility and value-at-risk. The data are large-cap stocks from the NYSE trade and quote database, for which the novel scheme's performance is found to be competitive.

C0503: Extreme events detection for volatility prediction

Presenter: **Andrea Montanino**, University of Naples Parthenope, Italy

Co-authors: Giovanni De Luca

The predictive power of extreme events such as financial bubbles and flash crashes on volatility is evaluated in cryptocurrencies and the Big Tech market. To precisely identify periodically collapsing bubbles and flash crashes, the backward supremum augmented Dickey-Fuller (BSADF) test is applied for date-stamping expansion and collapse phases. Volatility is then modeled using the most appropriate GARCH specification, with dummy variables included in the mean equation to capture each event type. The results show that these dummies significantly enhance predictive power not only for the reference asset but also cross-asset; for example, bubbles and crashes in Bitcoin anticipate volatility in other cryptocurrencies. Finally, a Diebold-Mariano test confirms that adding these dummies as external regressors yields a statistically significant forecasting improvement compared to the benchmark model without dummies. Understanding the link between extreme events and volatility is essential for both financial institutions and investors, as it informs risk management and portfolio allocation strategies.

C0519: Digital adoption and cyber security

Presenter: **Joann Jasiak**, York University, Canada

Co-authors: Pujee Tuvaandorj

The aim is to examine how Canadian firms balance the benefits of technology adoption against the rising risk of cybersecurity breaches. Data from the 2021 Canadian Survey of Digital Technology and Internet Use and the 2021 Canadian Survey of Cyber Security and Cybercrime are merged to investigate the trade-off firms face when adopting digital technologies to enhance productivity and efficiency, balanced against the potential increase in cyber security risk. The analysis explores the extent of digital technology adoption, differences across industries, the subsequent impacts on efficiency, and associated cybersecurity vulnerabilities. Aggregate variables are built, such as the business digital usage score and a cybersecurity incidence variable, to quantify each firm's digital engagement and cybersecurity risk. A survey-weight-adjusted Lasso estimator is employed, and a debiasing method for high-dimensional logit models is introduced to identify the drivers of technological efficiency and cyber risk. The analysis reveals a digital divide linked to firm size, industry, and workforce composition. While rapid expansion of tools such as cloud services or artificial intelligence can raise efficiency, it simultaneously heightens exposure to cyber threats, particularly among larger enterprises.

C1118: Well-conditioned covariance estimation via Bayesian eigenvalue regularization

Presenter: **Jesper Cremers**, Vrije Universiteit Brussel (VUB), Belgium

Co-authors: Kris Boudt, Steven Vanduffel, Kirill Dragun

Covariance matrices and their inverses are fundamental to a wide range of statistical applications. Traditional estimators often produce ill-conditioned or (nearly) singular matrices due to the presence of over-dispersed eigenvalues. We propose an explicit approach to regularize the eigenvalues of a class of covariance estimators using Bayes' theorem and a flexible prior that enforces positive semidefiniteness. A data-driven procedure determines the intensity of regularization, under which insights can be made into its behavior. The resulting estimators are both well-conditioned and accurate with respect to predictive performance. In a simulation study, we show that our proposed framework outperforms alternative rotation-invariant estimators. We validate its practical relevance through an empirical application to minimum-variance portfolios. Our approach significantly reduces out-of-sample variance and portfolio turnover.

CO255 Room BCB 213 HIGH-DIMENSIONAL AND STRUCTURAL APPROACHES TO SYSTEMIC RISK ASSESSMENT Chair: Leonardo Iania
C0484: Flow interdependence and systemic vulnerability in the U.S. mutual fund sector

Presenter: **Jie Yu**, UCLA Anderson School of Management, United States

Co-authors: Antoine Bouveret

Mutual funds can amplify financial shocks through asset liquidation driven by simultaneous redemptions. A novel copula-based framework is introduced to model the joint distribution of mutual fund flows, explicitly capturing nonlinear and tail-dependent relationships across fund categories. A new metric is proposed, *Flows in Distress (FiD)*, quantifying expected net outflows from a given fund category conditional on distress in another. Using U.S. mutual fund data spanning from 2000 to 2025, categories are identified that are particularly vulnerable to cross-fund spillovers and those acting as systemic transmitters during stress episodes. Additionally, the response of mutual fund flows is estimated to macro-financial shocks, demonstrating notable sensitivities to equity market volatility compared to Treasury market stress. Employing a Shapley value approach for systemic risk attribution, it is found that high-yield and world bond funds amplify systemic vulnerabilities, while investment-grade bond and hybrid funds generally mitigate redemption pressures, highlighting the nuanced roles different fund categories play during financial turmoil.

C0718: Large Bayesian structural matrix autoregression

Presenter: **Ignacio Moreira Lara**, Universitaet Duisburg Essen, Germany

Co-authors: Christoph Hanck, Jan Prueser

The structural identification of a high number of shocks in a vector autoregression (VAR) becomes increasingly challenging as the number of parameters proliferates and the set of restrictions needed for identification grows. These challenges hinder both the feasibility and interpretability

of spillover analysis in high-dimensional systems. A large structural Bayesian matrix autoregression is introduced that overcomes these obstacles by exploiting a Kronecker structured autoregressive structure, dramatically reducing both parameter and restriction counts. The algorithm is highly efficient and scales to large panels of time series. Contemporaneous relations are identified using established SVAR methods with minimal adjustments for the framework. To demonstrate its capabilities, the BMAR on monthly macro data is estimated for a broad panel of European countries, simultaneously extracting country-specific supply and demand shocks. The model delivers clear impulse responses and historical decompositions, providing concrete evidence of the asymmetric spillovers and spillbacks of multiple contemporaneous shocks across the euro-area.

C0865: Are oil price shocks priced in the cross-section of stock returns?

Presenter: **Thao Nguyen**, KBC Asset Management, Belgium

Co-authors: Leonardo Iania, Kristien Smedts, Liana Nersisyan

The purpose is to investigate the risk premia associated with oil-related shocks in the stock market. The shock components from a broad set of measures are extracted, including structural oil supply, demand, and inventory, market-based oil price expectations, news-based uncertainty indexes, and oil forecast disagreement. A negative risk premium associated with oil price shocks is documented: Stocks more sensitive to these shocks tend to earn lower expected returns. Furthermore, it is found that the greenest stocks are largely unaffected by oil price shocks, while browner stocks exhibit a negative relationship between oil shock exposure and expected performance.

C0871: The impact of oil price shocks in the corporate bond market

Presenter: **Liana Nersisyan**, UCLouvain, Belgium

Co-authors: Leonardo Iania, Marco Lyrio

The aim is to study the price of risk associated with oil price shocks on corporate bond risk premia. Oil price shocks are identified by using a structural vector autoregressive (SVAR) model of the global market for crude oil. These identified structural shocks - including oil supply, oil consumption demand, oil inventory demand, and global economic activity shocks - are then used as unspanned factors in a multi-market, single-pricing kernel framework. In the model, the first market represents a risk-free benchmark, while the second and third markets are represented by the yield curves on corporate bond indices of different rating classes (A and BBB). The analysis begins with predictive regressions of one-year holding period excess return of a corporate bond with different maturities on the first three principal components of zero-coupon yields of corporate bonds and oil-related macro factors. Preliminary findings showed that the impact of oil price shocks is stronger for lower-rated bonds, which suggests increased vulnerability in the firms with weaker credit profiles. Further empirical analysis is ongoing to deepen the understanding of the effects of the specific shocks and to refine the robustness of these results.

C0914: Forecasting lending rates in the euro area: Application of the advanced econometric methods

Presenter: **Vija Micune**, Latvia's Development Finance Institution Altum and University of Latvia, Latvia

The novel econometric approach of the quantile factor model is employed to deliver a more nuanced characterization and improved forecast of selected bank lending rates across multiple countries within the euro area. The research utilises, as explanatory variables in the model, not only traditional determinants of lending rates such as prevailing market interest rates, but also incorporates less conventional factors, including macro risk and credit risk. The accuracy and reliability of the results are evaluated through comparative analysis with alternative models and a combined density forecast. Findings support the initial hypothesis regarding the significance of the identified factors influencing lending rates across euro area countries. Although the pass-through from market rates to lending rates remains substantial during periods of positive interest rates, risk factors should not be underestimated as potential determinants of lending rates within the euro area. The application of the advanced econometric methods to the analysis and forecasting of lending rates across euro area countries yields more accurate predictions and deeper insight into the prospective evolution of the studied financial indicator.

C1025: Beyond financial markets: High-dimensional country risk assessment using latent factor models

Presenter: **Elias Wolf**, European Stability Mechanism, Luxembourg

Co-authors: Robert Blotvogel

In today's global economy, risks to financial stability are increasingly broad-based and not limited to developments in financial markets alone. Assessing country risk, therefore, requires a comprehensive approach that incorporates classic economic vulnerabilities across sectors, e.g., households, financial institutions, and corporates, as well as unexpected and external shocks that may arise from structural economic trends, political and institutional developments, or climate events and geopolitical tensions. Additionally, these risks might not be idiosyncratic but often interact in complex ways and propagate across time and space. Even countries with strong fundamentals may face financial stability risks if macroeconomic outcomes deviate from expectations due to adverse geopolitical shocks. Therefore, a new indicator is constructed, pooling information from a comprehensive high-dimensional set of individual risk sources to quantify and track multidimensional risks for different European countries. Following methodologies similar to the NFCI and CISS index, the strength of dynamic factor models is leveraged to handle big data sets with different sampling frequencies to identify latent risk factors. Simultaneously, it is investigated how the different sources contribute to the overall country risk over time. This approach aims to improve risk monitoring in the Euro area, where existing individual financial or economic indicators may underestimate emerging vulnerabilities.

CO194 Room BCB M202 NEW TOOLS FOR CAUSAL INFERENCE AND FORECASTING

Chair: Rustam Ibragimov

C0466: Efron-Stein type inequalities for randomly stopped processes

Presenter: **Victor H de la Pena**, Columbia University, United States

Efron-Stein type inequalities are established for randomly stopped processes, extending classical concentration results to settings with data-adaptive stopping times. The main result provides variance bounds for fixed functions of independent random variables evaluated at random stopping times that are adapted to the data but independent of the target function values. It carefully distinguishes between trivial cases (where stopping times are independent of the entire process) and genuinely non-trivial applications. The theoretical foundation relies on recent refinements of maximal inequalities for randomly stopped sums.

C0726: Score augmented Frobenius distance with applications in causal inference

Presenter: **Siyun He**, University of Michigan, United States

A novel method for causal inference is proposed in panel settings with network-valued outcomes by introducing the score augmented Frobenius distance, a metric that compares networks after sorting their adjacency matrices by node-level structural scores. These scores, which incorporate both observed covariates and structural features (e.g., centrality), serve to align nodes across networks under the assumption that structural roles, rather than identities, drive treatment effects. This sorting induces an equivalence class over node permutations, allowing valid comparisons between networks with unobserved heterogeneity. It is shown how the causal Frobenius distance can be used to extend standard difference-in-differences and synthetic control methods to settings where the outcome is a network. The framework is applicable to a variety of empirical settings, including social networks, trade networks, and institutional relationships, where interventions affect structural properties rather than specific node labels. Formal identification results are provided, asymptotic behavior is discussed under network sampling, and performance on simulated and real-world policy interventions is demonstrated.

C1051: Implied event risk from option prices

Presenter: **Johan Walden**, UC Berkeley, Haas School of Business, United States

Co-authors: Richard Stanton, Eben Lazarus

A nonparametric method is introduced for inferring risk distributions of anticipated events of publicly traded firms, such as earnings announcements, securities litigation events, and the resolution of takeover bids and proxy fights. The method is applied to a sample of specific events, including earnings announcements, product launches, and FOMC announcements, and it is shown how well-known characteristics of asset prices before such events, e.g., concave implied volatility curves, are explained by anticipated event risk. The method provides new insight on the effect of anticipated event risk on asset prices, how to measure such risk, and — more specifically — a fruitful way of adjusting option prices to account for such risk close to maturity dates.

C1224: Market timing with bi-objective cost-sensitive machine learning

Presenter: **Artem Prokhorov**, University of Sydney, Australia

The aim is to develop a framework for cost-sensitive training of machine learning models that predict the direction of aggregate stock returns. A bi-objective loss function is designed that augments the traditional log-loss objective with an objective that minimizes the cost of individual false-positive and false-negative classification errors. It is argued that the option-implied conditional value-at-risk is a natural measure of the misclassification costs in such models. The bi-objective optimization framework permits us to isolate the effect of cost-sensitivity from log-loss minimization, and to integrate forward-looking information from options markets directly into the model training process. Changes are studied in the classification performance of elastic-net logistic regression and gradient-boosted decision trees trained using the bi-objective framework. The new approach improves the risk-adjusted returns of market timing strategies and substantially reduces downside risk.

C1269: Robust Cauchy-based methods for predictive regressions

Presenter: **Rustam Ibragimov**, Imperial College Business School and New Economic School, United Kingdom

Co-authors: Jihyun Kim, Anton Skrobotov

The purpose is to develop robust inference methods for predictive regressions, addressing challenges posed by endogenously persistent or heavy-tailed regressors, as well as persistent volatility in the errors. Building upon the Cauchy estimation framework, two novel tests are proposed: One that relies on t-statistic-based group inference and another that employs a hybrid approach combining Cauchy and OLS estimation. These methods effectively tackle key issues in standard inference procedures, including size distortions arising from endogenously persistent or heavy-tailed regressors and persistent volatility dynamics. The proposed methods are straightforward to implement and broadly applicable to both continuous and discrete time models. Extensive simulation studies highlight the finite-sample advantages of the proposed methods under realistic settings. An empirical application is provided to test the predictability of excess returns for two major stock indices using the dividend-price and earnings-price ratios as predictors. The results indicate that the dividend-price ratio has predictive power, while the earnings-price ratio does not significantly predict returns.

CO100 Room BCB 308 RESEARCHER DEGREES OF FREEDOM, FLEXIBLE MODELING AND INTERPRETABILITY Chair: Roman Hornung

C0821: Comparing methods for flexible estimation of non-linear associations: The importance of suitable performance measures

Presenter: **Theresa Ullmann**, Medical University of Vienna, Austria

Flexible regression techniques (e.g., spline-based approaches or fractional polynomials) allow for modeling non-linear associations between continuous predictors and outcomes, often leading to improved model performance compared to specifying the associations as linear. Simulation studies are a key tool for comparing such techniques, but their conclusions critically depend on the performance measures used to evaluate how well the estimated curves recover the true underlying functions. A systematic categorization of performance measures is presented for evaluating estimated non-linear associations between an outcome and continuous predictors. To illustrate the practical implications of these choices, examples are presented that highlight how different performance measures can favor different methods. Results are also shown from a simulation study comparing several flexible modeling approaches, using performance measures from the proposed categorization. The emphasis is on the importance of aligning performance measures with the aim and the encouragement of more transparent and thoughtful evaluation strategies in methodological research.

C0938: Researcher degrees of freedom and associated challenges in designing statistical parametric simulation studies

Presenter: **Felix Julian David Lange**, LMU Munich, Germany

Parametric simulation studies are indispensable for evaluating and comparing statistical methods. Designing such a study involves numerous decisions, and researchers generally have a great deal of flexibility when making decisions about the various design aspects, including parameter values, performance measures, and methods involved. While this flexibility, commonly referred to as researcher degrees of freedom, is part of the appeal of simulation studies, it also presents considerable challenges and risks. Not only is there a risk that the large number of decisions leads to arbitrary choices. The researcher-specified conditions might also be unrealistic, which is problematic because statistical simulation studies are frequently used as a basis for practical recommendations about methods. Researchers may also unintentionally or intentionally make biased choices that favor certain methods or outcomes, such as demonstrating that a particular method is superior. These issues, along with some potential remedies, are discussed. Additionally, it is illustrated how the researcher's degrees of freedom in the design can affect a study's results and lead to findings that do not generalize well.

C0920: Exploring meaningful analytical choices using the vibration of effects framework

Presenter: **Simon Lemster**, LMU Munich, Germany

In many studies in biomedical research, contradictory results are found even when investigating the same research question. Effects often vary considerably across studies, not only due to sampling variability but also due to the multiplicity of possible analysis strategies. These include the exact definition of exposure and outcome, the selection of adjustment variables, or data preprocessing decisions, such as treatment of outliers and missing data. The influence of these analysis decisions can be investigated with the vibration of effects framework. Using the frequently studied example of body composition and cardiovascular mortality, this approach is demonstrated on data from a German population-based cohort study. Findings show that while a positive association between body composition and cardiovascular disease often emerges, the strength and direction of the effect are highly sensitive to analytic decisions. In some model specifications, the association may even reverse. However, not all specifications are equally plausible from a causal or methodological perspective. Some choices may introduce bias, for example, through adjustment for colliders or inappropriate data handling. Therefore, restricting the analysis to a subset of theoretically justifiable decisions may reduce the observed vibration of effects. This could increase both credibility and practical utility of multiverse analyses in applied statistical research.

C0690: Permutation-based multiple testing-controlled variable selection using random forests

Presenter: **Tim Mueller**, Staburo GmbH, Germany

Co-authors: Roman Hornung, Silke Szymczak, Hannes Buchner

Identifying relevant biomarkers is critical in clinical research and precision medicine, particularly when analyzing high-dimensional data. Random forests (RFs) are promising for such settings due to their flexibility, ease of use, and their ability to handle datasets with more variables than samples. RFs assess the importance of each variable in predicting the outcome using variable importance (VIMP) scores. However, the lack of a known statistical distribution of VIMP scores prevents standard statistical testing and associated multiple testing adjustment for the purpose of variable selection. The aim is to propose a novel method for multiple testing-controlled variable selection. The approach, similar to permutation testing, involves generating permuted counterparts for each variable and comparing their VIMPs across iterations to calculate p-values. However,

unlike competing methods, the correlation structure is preserved between the covariates in the permutations to guard against biases. With promising results, the method is evaluated against three competing RF variable selection approaches in simulations that involve high- and low-dimensional data, as well as correlated and categorical variables. Moreover, it is applied to a real dataset to demonstrate its practical use. The method's results integrate seamlessly into standard VIMP plots, providing a flexible and transparent way to interpret results in a familiar format.

C0833: Unity forests: Improving interaction modelling and interpretability in random forests

Presenter: **Roman Hornung**, University of Munich, Germany

Random forests (RFs) are widely used for prediction and variable importance analysis and are often assumed to capture interactions via recursive splitting. However, since the splits are chosen locally at each node, RFs capture interactions only when at least one involved covariate has a notable marginal effect. Unity forests (UFOs) are introduced, an RF variant designed to exploit interactions involving covariates without marginal effects. In UFOs, the first few splits of each tree (by default three) are optimized jointly across a random covariate subset to form a "tree root" that better captures such interactions; the remainder is grown conventionally. A new variable importance measure (VIM) is also proposed - the unity VIM - based on the out-of-bag split criterion values from the tree roots. Only a small fraction of root splits with the highest in-bag values are considered per covariate, reflecting that covariates involved in interactions tend to be influential in relatively few trees. In a simulation study, the unity VIM consistently ranked interacting covariates without marginal effects above non-influential ones - unlike conventional RF-based VIMs. In a large-scale real data-based comparison, UFOs improved predictive accuracy and discrimination over standard RFs, with similar calibration. Finally, selecting representative trees is explored per covariate from the tree roots, offering interpretable insight into each covariate's individual or interactive predictive role.

CO095 Room BCB 309 TOPICS IN STATISTICAL GENETICS AND BIOINFORMATICS

Chair: Somak Dutta

C0264: A hybrid mixture approach for clustering and characterizing cancer data

Presenter: **Fan Dai**, Michigan Technological University, United States

Co-authors: Kazeem Kareem

Model-based clustering is widely used for identifying and distinguishing types of diseases. However, modern biomedical data, coming with high dimensions, make it challenging to perform the model estimation in traditional cluster analysis. The incorporation of factor analyzer into the mixture model provides a way to characterize the large set of data features, but the current estimation method is computationally impractical for massive data due to the intrinsic slow convergence of the embedded algorithms, and the incapability to vary the size of the factor analyzers, preventing the implementation of a generalized mixture of factor analyzers and further characterization of the data clusters. A hybrid matrix-free computational scheme is proposed to efficiently estimate the clusters and model parameters based on a Gaussian mixture, along with generalized factor analyzers to summarize the large number of variables using a small set of underlying factors. The approach outperforms the existing method with faster convergence while maintaining high clustering accuracy. The algorithms are applied to accurately identify and distinguish breast cancer based on large tumor samples, and to provide a generalized characterization for subtypes of lymphoma using massive gene records.

C0535: R/Shiny web applications for making gene expression data accessible to non-bioinformaticians

Presenter: **Garrett Dancik**, Eastern Connecticut State University, United States

Repositories such as the Gene Expression Omnibus (GEO) and the Genomic Data Commons contain gene expression profiles for millions of samples. Yet, the retrieval and analysis of this data is difficult without bioinformatics and statistical expertise. Several web applications are described to facilitate the analysis of gene expression data. ShinyGEO is a web application that allows users to download gene expression data sets directly from GEO, select a gene of interest, and carry out differential expression and survival analyses with customizable graphics. The Bladder Cancer Biomarker Evaluation Tool (BC-BET) allows for rapid evaluation of candidate diagnostic and prognostic gene expression biomarkers in 15 bladder cancer patient cohorts (N = 1559). Ongoing development of a similar Acute Myeloid Leukemia Biomarker Evaluation Tool will also be described, with an emphasis on the importance of data processing and analysis using robust statistical methods. The availability of these tools makes gene expression data accessible to non-bioinformaticians, promising to lead to a better understanding of biological processes and genetic diseases such as cancer.

C0955: A branching process model for digital read quantification

Presenter: **Karin Dorman**, Iowa State University, United States

Sequenced read counts are a ubiquitous data summary of modern high-throughput biological methods used to observe metagenomes, genomes, transcriptomes, epigenomes, and various kinds of molecular interactions and functions. Almost all such count data are obtained after amplification of sampled molecules, which can bias and overdisperse the biological signal of interest. The purpose is to develop and investigate a novel model for count data that better adheres to the experimental generative process than Poisson and negative binomial models with or without zero-inflation. The model is based on a branching process model of polymerase chain reaction (PCR) amplification. It naturally accounts for overdispersion and zero inflation, with meaningful parameters directly linked to biological processes. In particular, the first estimates of PCR amplification efficiency are provided during library preparation and estimate the effects of primer mismatch on sampling efficiency.

C1005: The causality and dynamics of gene regulation during macrophage-neutrophil differentiation

Presenter: **Manu Manu**, University of North Dakota, United States

Co-authors: Yen Lee Loh, Trevor Long, Nimasha Samarawickrama

During gene regulation, DNA accessibility is thought to limit the availability of transcription factor (TF) binding sites, while TFs can increase DNA accessibility to recruit additional factors that upregulate gene expression. Given this interplay, the causative regulatory events in the modulation of gene expression remain unknown for the vast majority of genes. The causality of gene regulation is investigated during macrophage-neutrophil differentiation using joint RNA-Seq and ATAC-Seq time series datasets. It is shown using an unsupervised learning technique that only 10 distinct temporal expression patterns are sufficient to recapitulate the expression of 36,000 transcripts with high fidelity. Furthermore, information transfer during differentiation was found to occur in cascading waves of gene expression culminating in the permanent turning on of certain genes after 80h. The genes enriched in each wave suggest a characteristic order of physiological remodeling-signal transduction, translation and mRNA processing, metabolism, and, ultimately, myeloid phenotypic processes. It is also found that gene expression dynamics do not observe any strict relationship with the dynamics of the accessibility of enhancers or promoters, varying from highly positively correlated to uncorrelated, to highly negatively correlated. In a detailed analysis of a key neutrophil gene, Cebpa, it is shown that TF occupancy rather than DNA accessibility is the causal driver of gene expression.

C1464: Deconvolution of cell free DNA via Bayesian tree-based marker selection and signature estimation

Presenter: **Christopher McKennan**, University of Pittsburgh, United States

Co-authors: Yucheng Wang

Cell-free DNA (cfDNA) consists of small DNA fragments released into the bloodstream during cell death and has enabled non-invasive diagnostics for many diseases. A key step in cfDNA analysis is cellular deconvolution, which estimates the proportion of cfDNA fragments originating from each major cell type using CpG methylation patterns. Existing methods suffer from major limitations, including poorly curated marker CpGs, inaccurate cell type-specific (CTS) methylation signatures, and slow runtimes. To address these challenges, cf-TREBLE is developed, a statistically rigorous and scalable method for identifying marker CpGs, estimating CTS methylation signatures, and performing deconvolution. cf-TREBLE uses CTS reference data, a hierarchical cell type tree, and a novel Bayesian model to classify CpGs based on whether their methylation

is shared across or unique to specific cell types. This enables the identification of highly reliable markers and improves CTS signature estimates by pooling information across related cell types, especially when reference data are sparse. cf-TREBLE then applies a novel, computationally efficient deconvolution algorithm that accounts for inter-subject variability in methylation. cf-TREBLEs superior performance is demonstrated on realistic simulated data, and it is applied to develop risk prediction models for adenomyosis and endometriosis.

CO065 Room BCB 310 NEW DEVELOPMENTS IN FUNCTIONAL DATA ANALYSIS
Chair: Eleonora Arnone
C0665: Inference on functional data through e-values

Presenter: **Alessia Pini**, Università Cattolica del Sacro Cuore, Italy

Co-authors: Simone Vantini

A very recent area of statistical inference proposes to replace p-values with e-values for testing hypotheses on univariate or multivariate data. Namely, an e-value (also known as a betting score) is the value taken by a random variable whose expected value is equal to one under the null hypothesis. We propose the extension of e-values to functional data. In detail, we start by defining a point-wise e-value function for performing inference on each point of the domain of the functions, separately. Then, we show that the integral average of the point-wise e-value function is a valid e-value for performing global inference on functional data, and discuss its properties and interpretation.

C0884: On identifying functional motifs and using them to aid forecasting and missing value imputation

Presenter: **Jacopo Di Iorio**, Emory University, United States

Co-authors: Francesca Chiaromonte, Marzia Cremona

Functional data analysis faces a novel challenge: Identifying functional motifs, or shapes, that may be repeated multiple times within each functional observation or across multiple curves belonging to the same set. To address this issue, funBAlign is introduced, a multi-step approach that employs agglomerative hierarchical clustering with complete linkage and functional distances based on mean squared residue scores and virtual error. These distances enable funBAlign to detect functional motifs that may be shifted or scaled along the y-axis. To validate the effectiveness of the methodology, simulations and case studies that demonstrate its ability to identify functional motifs are presented. The identification of functional motifs can be leveraged to solve other important problems in functional data analysis. For instance, portions characterized by the same motif are hypothesized to evolve similarly, which can aid in forecasting and the missing portions imputation problems.

C0941: Differentially private geodesic regression

Presenter: **Carlos Soto**, University of Massachusetts Amherst, United States

Co-authors: Aditya Kulkarni

In statistical applications, it has become increasingly common to encounter data structures that live on nonlinear spaces such as manifolds. Classical linear regression, one of the most fundamental methodologies of statistical learning, captures the relationship between an independent variable and a response variable, which both are assumed to live in Euclidean space. Thus, geodesic regression emerged as an extension where the response variable lives on a Riemannian manifold. The parameters of geodesic regression, as with linear regression, capture the relationship of sensitive data, and, hence, one should consider the privacy protection practices of said parameters. Releasing differentially private (DP) parameters of geodesic regression is considered via the k-norm gradient (KNG) mechanism for Riemannian manifolds. Theoretical bounds are derived for the sensitivity of the parameters, showing they are tied to their respective Jacobi fields and hence the curvature of the space. This corroborates recent findings of differential privacy for the Fréchet mean. The efficacy of the methodology is demonstrated on the sphere, $\mathbb{S}^d \subset \mathbb{R}^{d+1}$, and since it is general to Riemannian manifolds, the manifold of Euclidean space, which simplifies geodesic regression to a case of linear regression. The methodology is general to any Riemannian manifold and thus it is suitable for data in domains such as medical imaging and computer vision.

C1120: Causal inference meets functional data: On the estimation of functional average and conditional treatment effects

Presenter: **Lorenzo Testa**, Carnegie Mellon University, United States

Co-authors: Francesca Chiaromonte, Edward Kennedy, Tobia Boschi, Filippo Salmaso

Understanding causal relationships is a central goal in science, but traditional methods often fall short when dealing with complex data. This challenge is pronounced in applications where outcomes are not simple scalars but are instead functions observed over a continuous domain. The purpose is to introduce a unified framework for causal inference with functional outcomes. The problem of estimating the Functional Average Treatment Effect (FATE) is addressed to understand the overall impact of an intervention across a population. The aim is to present a novel, doubly robust estimator that guarantees consistent estimation even if one of the underlying models – either for treatment assignment or outcome – is misspecified. Its theoretical properties are established, ensuring valid inference through the construction of simultaneous confidence bands. Building on this foundation, the more nuanced challenge of personalization is then tackled by estimating the functional conditional average treatment effect (F-CATE). A novel meta-learning framework is introduced, designed to uncover how functional treatment effects vary across individuals. This approach is also doubly robust, integrating advanced functional regression techniques to provide reliable, individualized causal insights. Across both problems, the robustness of the methods is demonstrated through simulations. Their real-world utility in uncovering meaningful causal effects from complex health data is illustrated.

C1346: Finite-sample improvements in nonparametric functional regression through weighted pseudo-metrics

Presenter: **Kwo Lik Lax Chan**, Università degli Studi del Piemonte Orientale, Italy

Co-authors: Laurent Delsol, Aldo Goia

One of the main problems in functional data analysis is the selection of pseudo-metric as a distance measure between curves; in particular, in the nonparametric regression context, it has a direct impact as it captures the information contained in the explanatory curve by extracting the informative features of the explanatory curve. The idea of weighted pseudo-metric is introduced, implemented, and discussed. Performances of the choices of pseudo-metric used are evaluated by means of a Monte Carlo Study.

CO264 Room BCB 311 STATISTICAL METHODS FOR ECOLOGICAL AND ENVIRONMENTAL COMPLEXITY
Chair: Crescenza Calculi
C0567: Quantifying the unexpected: Developing statistical tests for spatial entropy in ecology and landscape analysis

Presenter: **Linda Altieri**, University of Bologna, Italy

Shannon entropy is widely used in ecology, biodiversity, and landscape studies to quantify the heterogeneity of a system. Its appeal lies in its versatility and broad applicability, including to qualitative and categorical data. However, its known limitations in capturing spatial structure have led to the development of alternative entropy-based measures tailored for spatial data, such as those proposed by Batty, Karlstrom, O'Neill, Leibovici, and Altieri, among others. While these spatial entropy measures hold great potential for describing complex spatial patterns, they ultimately reduce system complexity to a single value. As such, it remains difficult to assess whether the observed heterogeneity meaningfully deviates from an expected level, either under a null model or in comparison to other configurations. To date, statistical testing procedures for entropy values have only been developed for Shannon entropy in non-spatial settings, leaving a gap in the spatial domain. The development of formal statistical tests for spatial entropy measures is proposed. These tests aim to provide a framework for assessing the significance of observed spatial heterogeneity, supporting more rigorous interpretation and comparison across spatial systems.

C0382: A cross-validation framework for log-Gaussian Cox processes with presence-only data in marine ecology

Presenter: **Daniele Poggio**, Politecnico di Torino, Italy

Co-authors: Gian Mario Sangiovanni, Gianluca Mastrantonio, Giovanna Jona Lasinio, Daniele Ventura, Edoardo Casoli, Stefano Moro

Building on a recent study, a novel implementation of a novel cross-validation strategy is proposed for model selection within the framework of log-Gaussian Cox processes (LGCP), designed to enhance predictive performance assessment. The approach implements a k-fold cross-validation scheme based on a spatial thinning mechanism to define the folds. Model comparison is performed by computing posterior distributions of raw residuals across a partitioned spatial domain and summarizing them using the continuous ranked probability score (CRPS), yielding a single metric per model. The framework is computationally efficient and suitable for integrating heterogeneous presence-only data, with particular relevance to marine ecology. Its application is demonstrated in modeling the spatial distribution of Holothurians recorded during nine survey campaigns (2022/2024) around Giglio Island, Italy. The surveys integrated high-resolution photogrammetry with diver-based visual censuses, resulting in variable detection probabilities across habitats, especially in Posidonia Oceanica meadows. To address this complexity, LGCP models are adopted with a shared spatial Gaussian process component, capturing habitat structure, environmental covariates, and temporal variability. Multiple model specifications are evaluated, each incorporating different sets of spatial covariates reflecting habitat and geomorphological features, and the most predictive model is identified.

C0440: From object detection to modeling spatial intensity of marine biodiversity

Presenter: **Gian Mario Sangiovanni**, Sapienza University, Italy

Co-authors: Daniele Poggio, Gianluca Mastrantonio, Giovanna Jona Lasinio, Daniele Ventura, Stefano Moro

In ecology, photogrammetry is a crucial method for efficiently acquiring non-destructive samples of natural environments. When the goal is to estimate the spatial distribution of animals, detecting objects in large-scale images becomes essential. Object detection models enable large-scale analysis but introduce uncertainty, as the probability of detection depends on various factors. A key aspect of this process is the selection of the confidence threshold used during detection. A conservative threshold ensures high precision but reduces sensitivity, which can lead to an underestimation of community size and bias in species distribution models. The purpose is to utilize YOLOv11; however, the main advantage of the approach is its flexibility, allowing the usage of any detector. To address detection bias, the distribution of holothurians (sea cucumbers) is modeled in an area near the coast of Giglio Island using a thinned non-homogeneous Poisson process (NHPP). It assumes that a "true" intensity function accurately describes the distribution, while the observed process, resulting from independent thinning, is represented by a "degraded" intensity. The detection function regulates the thinning mechanism, influenced by the object's location and other detection-related features.

C0886: Comparing spatiotemporal temperature models for heat-related mortality risk assessment in Lazio, Italy

Presenter: **Emiliano Ceccarelli**, Sapienza University of Rome, Italy

Co-authors: Giovanna Jona Lasinio, Giada Minelli, Marta Blangiardo, Jorge Castillo-Mateo, Sandra Gudziunaite

The aim is to investigate how different spatiotemporal temperature models affect the estimation of heat-related mortality in Lazio, Italy (2008-2022). Three methods are compared to reconstruct daily maximum temperature at the municipality level: Two Bayesian station-based approaches, a quantile autoregressive model with spatial interpolation, a Gaussian model via INLA-SPDE, and a satellite-based method using ERA5. Station-based models show higher and more spatially variable temperatures than satellite-based ones, especially in warmer provinces. Using individual mortality data for cardiovascular and respiratory causes, temperature-mortality associations are estimated through Bayesian conditional Poisson models in a case-crossover design. Exposure is defined as the mean maximum temperature over the previous three days. Additional models include heatwave definitions combining different thresholds and durations. All models show a marked increase in relative risk at high temperatures, but the temperature of minimum risk varies notably across methods. Station-based models estimate higher minimum-risk temperatures compared to ERA5. Stratified analyses reveal higher RR increases in females and the elderly (80+). Heatwave effects depend on definitions, but all methods capture the prolonged heat exposure effect. Overall, results confirm the importance of temperature model choice in epidemiology and provide insights for early warning systems and climate-health adaptation strategies.

C0662: Machine learning models for predicting the socio-economic impacts of climate change

Presenter: **Angela Maria D'Ugento**, University of Bari Aldo Moro, Italy

Co-authors: Margaret Antonicelli

Climate change has caused enormous environmental damage through forest fires and floods. These extreme events also have social, political, and economic consequences. To understand these delayed effects of climate change, it is important to consider some phenomena as environmental problems. Climate justice focuses on climate change as an ethical and political issue by highlighting how the most vulnerable populations, despite contributing the least to the problem, suffer the most from its impacts, making it a human rights issue. Climate migration is largely determined by the availability and accessibility of resources and exposure to environmental hazards. Gender inequality could also be exacerbated by climate change, as women are more vulnerable due to the division of labor, restrictive gender norms, and underrepresentation in climate-related decision-making. Finally, economic impacts such as the rise in coffee and cocoa prices also have a direct effect on people's lives. To analyze climate change trends, an analysis of anomalous temperature data in European countries from 1850 to 2023 is carried out using machine learning prediction models. Conventional linear models have struggled to accurately capture climate trends, even with SARIMA. A long short-term memory (LSTM) neural network optimized with a modified Adam algorithm showed superior prediction performance.

CO327 Room BCB 312 STATISTICAL AI IDEAS FOR MODERN BIOSTATISTICS

Chair: Marcos Matabuena

C0477: Denoising data with measurement error using a reproducing kernel-based diffusion model

Presenter: **Ruoyu Wang**, Harvard School of Public Health, United States

Co-authors: Mingyang Yi, Marcos Matabuena

The ongoing technological revolution in measurement systems enables the acquisition of high-resolution samples in fields such as engineering, biology, and medicine. However, these observations are often subject to errors from measurement devices. Motivated by this challenge, a denoising framework that employs diffusion models is proposed to generate denoised data whose distribution closely approximates the unobservable, error-free data, thereby permitting standard data analysis based on the denoised data. The key element of the framework is a novel reproducing kernel Hilbert space-based method that trains the diffusion model with only error-contaminated data, admits a closed-form solution, and achieves a fast convergence rate in terms of estimation error. Furthermore, the effectiveness of the method is verified by deriving an upper bound on the Kullback-Leibler divergence between the distributions of the generated denoised data and the error-free data. A series of conducted simulations also verifies the promising empirical performance of the proposed method compared to other state-of-the-art methods. To further illustrate the potential of this denoising framework in a real-world application, it is applied in a digital health context, showing how measurement error in continuous glucose monitors can influence conclusions drawn from a clinical trial on diabetes Mellitus.

C0531: Conformal alignment: Knowing when to trust foundation models with guarantees

Presenter: **Ying Jin**, University of Pennsylvania, United States

Before deploying outputs from foundation models in high-stakes tasks, it is imperative to ensure that they align with human values. For instance, in radiology report generation, reports generated by a vision-language model must align with human evaluations before their use in medical decision-making. Conformal alignment is presented, a general framework for identifying units whose outputs meet a user-specified alignment criterion. It is guaranteed that on average, a prescribed fraction of selected units indeed meet the alignment criterion, regardless of the foundation model or the data distribution. Given any pre-trained model and new units with model-generated outputs, conformal alignment leverages a set of reference data with ground-truth alignment status to train an alignment predictor. It then selects new units whose predicted alignment scores surpass a

data-dependent threshold, certifying their corresponding outputs as trustworthy. Through applications to question answering and radiology report generation, it is demonstrated that the method is able to accurately identify units with trustworthy outputs via lightweight training over a moderate amount of reference data. En route, the informativeness of various features is investigated in alignment prediction, and they are combined with standard models to construct the alignment predictor.

C0749: Preference-integrated dynamic treatment regimes: Methodology and application to SMARTs

Presenter: **Yating Zou**, University of North Carolina at Chapel Hill, United States

Co-authors: Michael Kosorok, Lesile Wilson, Joshua Zitovsky

In modern healthcare, it is often useful for treatment decisions to balance multiple outcomes according to patient preferences. For example, efficacy and side effects. The purpose is to introduce latent utility Q-learning (LUQ-Learning), a Q-learning algorithm augmented by latent variable modeling, that embeds random individual preferences into the optimality criterion. LUQ-Learning supports a finite number of decision points and finite dimensions of outcomes with asymptotic guarantees under realistic assumptions. In simulations, LUQ-Learning consistently outperforms alternatives. The method was applied to BEST, a SMART study on chronic low back pain, highlighting its practical utility in precision medicine.

C0851: Nuisance parameter tuning for inference in observational studies

Presenter: **Rajarshi Mukherjee**, Harvard T.H. Chan School of Public Health, United States

The purpose is to discuss the issue of nuisance parameter tuning for estimating quantities in observational studies, such as the average treatment effect and measures of conditional dependence. Typical methods of estimating such quantities of interest rely on estimating nuisance functions often through the lens of nonparametric and/or high-dimensional machine learning methods. Whereas many popular ideas pertain to tuning these nuisance function estimation from a prediction perspective and subsequently perform downstream bias correction for valid inference of low dimensional summaries of interest in the observational studies of interest, cases are explored to show that there exists a delicate interplay between nuisance function estimation strategies, type of estimators that uses these nuisance functions in its pipeline of estimation of the final object of interest, and sample splitting strategies that are now popular to allow flexible methods of nuisance function estimation without jeopardizing the standard errors of estimators of the downstream objects of interest. The above is explored through the lens of specific functionals that arise in the context of causal inference, and both are studied in nonparametric and high-dimensional regimes.

C0924: Sparse PCA with multiple components

Presenter: **Jean Pauphilet**, London Business School, United Kingdom

Co-authors: Ryan Cory-Wright

Sparse principal component analysis (sPCA) is a cardinal technique for obtaining combinations of features, or principal components (PCs), that explain the variance of high-dimensional datasets in an interpretable manner. This involves finding directions (or PCs) that are jointly sparse and mutually orthogonal. Most existing works address sparse PCA (via methods such as iteratively computing one sparse PC and deflating the covariance matrix) that do not guarantee the orthogonality, let alone the optimality, of the resulting solution. Several optimization approaches are presented (semidefinite, combinatorial, and non-convex) to derive both upper bounds (on the maximum amount of variance explainable) and feasible solutions (i.e., sparse orthogonal PCs). Numerically, the algorithms match (and sometimes surpass) the best performing methods in terms of fraction of variance explained and systematically return PCs that are sparse and orthogonal. In contrast, it is found that existing methods like deflation return solutions that violate the orthogonality constraints, even when the data is generated according to sparse orthogonal PCs.

C0212 Room BCB 313 DEVELOPMENTS IN RECEIVER OPERATING CHARACTERISTIC CURVE ANALYSIS

Chair: Andrew Spieker

C0760: Semi-parametric confidence bands for the ROC curve based on unions of transformed elliptical regions

Presenter: **Andrew Spieker**, Vanderbilt University Medical Center, United States

Receiver operating characteristic (ROC) analysis is a helpful tool to evaluate the capacity of a measure to distinguish between states. Assumptions regarding the structure of the ROC curve can improve the efficiency of estimation for the ROC curve itself or its associated area under the curve (AUC). Many are familiar with nonparametric methods for generating an ROC curve; fewer are familiar with parametric and semi-parametric methods. Because the ROC curve is a two-dimensional projection of a fundamentally three-dimensional object, the challenges of characterizing uncertainty in the ROC curve or its associated AUC are often overlooked. The focus is on methods to characterize uncertainty based on a unions-of-transformed-ellipses approach.

C0218: A comprehensive overview of ROC curve estimation: Applications in medical research

Presenter: **Musie Ghebremichael**, Harvard University, United States

The receiver-operating characteristic (ROC) curve is widely used in many fields, including medicine. Due to its widespread usefulness, various ROC estimation and testing methods have been continuously evolving. Each method makes different assumptions about the underlying distribution of the data, yields a curve with different properties, and produces different estimates of the subsequent area under the ROC curve. These often leave researchers with a complex problem of deciding which method to use. The properties and performances of parametric, nonparametric, semiparametric, Bayesian, and placement value-based ROC curves are investigated. Extensive simulation studies were carried out across various distribution types, sample sizes, and degrees of overlap between the diseased and non-diseased population distributions. Additionally, real-world data from pediatric HIV/AIDS research were analyzed to illustrate and compare the methods. Findings demonstrate that different ROC estimation methods can produce significantly varying results, highlighting the necessity for thorough validation of chosen approaches. With the growing focus on biomarker discovery and the advent of high-throughput assays and diagnostic tools, the ROC will remain an important analytic tool in medical research. A practical guide is provided for statisticians and researchers working in medical fields to select ROC inference methods appropriate to their specific data.

C0223: Receiver operating characteristic curve for complex survey data

Presenter: **Tamy Tsujimoto**, Google - YouTube, United States

Co-authors: Jianwen Cai

The receiver operating characteristic (ROC) curve is frequently used to evaluate the accuracy of medical diagnostic tests. Currently, analysis based on the ROC curve has been performed in large public-use data arising from complex survey samples by ignoring the sampling scheme. A nonparametric estimator is proposed for the ROC curve that accounts for complex survey sampling. The asymptotic properties of the estimator are developed using empirical process arguments. Simulation studies showed that the proposed estimator performed well in the practical situations considered, with better performance for larger sample size and disease proportions. The estimator was illustrated in the National Health and Nutrition Examination Survey (NHANES) to evaluate the discrimination of a traditional risk calculator for undiagnosed diabetes

C0357: Nonparametric worst-case bounds for publication bias on the summary receiver operating characteristic curve

Presenter: **Satoshi Hattori**, The University of Osaka, Japan

The summary receiver operating characteristic (SROC) curve has been recommended as one important meta-analytical summary to represent the accuracy of a diagnostic test in the presence of heterogeneous cutoff values. However, the selective publication of diagnostic studies for meta-analysis can induce publication bias (PB) in the estimate of the SROC curve. Several sensitivity analysis methods have been developed to quantify PB on the S ROC curve, and all these methods utilize parametric selection functions to model the selective publication mechanism. A new

sensitivity analysis approach is proposed that derives the worst-case bounds for the SROC curve by adopting nonparametric selection functions under minimal assumptions. The estimation procedures of the worst-case bounds use the Monte Carlo method to approximate the bias on the SROC curves along with the corresponding area under the curves, and then the maximum and minimum values of PB under a range of marginal selection probabilities are optimized by nonlinear programming. The proposed method is applied to real-world meta-analyses to show that the worst-case bounds of the SROC curves can provide useful insights for discussing the robustness of meta-analytical findings on diagnostic test accuracy.

C0770: Comparison of diagnostic tests in multireader multicase ROC studies with missing data: A rank-based approach

Presenter: **Guangyong Zou**, Western University, Canada

Comparison of diagnostic tests can be done by quantifying differences between areas under the receiver operating characteristic curves (AUCs) using the multireader multicase (MRMC) design, where each case undergoes each of several diagnostic tests and the resulting images are interpreted by each of several readers who are blinded to the true disease status. This design is required for the clinical evaluation of computer-aided diagnostic devices and imaging diagnostic modalities by regulatory agencies. Procedures for confidence interval estimation are under-developed for MRMC, especially in the presence of missing data. We propose a rank-based approach to confidence interval estimation for AUCs and their difference. We first transform each observation into a case-specific accuracy value based on the difference between its rank among all observations in the reader-test combination and its rank within its disease status group. We then analyze these case-specific accuracies using regression models for clustered data to obtain point and variance-covariance estimates of the AUCs. Asymptotic confidence intervals for test-specific AUCs and their difference are constructed under logit- and inverse hyperbolic tangent (artanh)-transformation, respectively. Simulation results based on real studies suggest that our method performed very well in terms of coverage and empirical power, even in the presence of missing data.

CO110 Room BCB 402 SELECTED WORK FROM CAUCUS FOR WOMEN IN STATISTICS (CWS) MEMBERS

Chair: Vanda Lourenco

C0378: Contrasting and interpreting test decisions induced by information borrowing in hybrid-control clinical trial designs

Presenter: **Silvia Calderazzo**, German Cancer Research Center (DKFZ), Germany

When designing a novel clinical trial, some external information about the control and/or treatment effect is typically available. Borrowing of such external information is often desired in order to improve the trials' efficiency, and Bayesian designs with informative prior distributions offer a natural framework to achieve this aim. One major concern in this context is the potential for heterogeneity between the current and external data sources, which can significantly increase both the chances of making erroneous test decisions and estimation error. Robust dynamic prior choices are thus often employed: Such priors gradually discount external information based on the observed heterogeneity between the two information sources. In this context, it is of interest to understand how external information borrowing in general, and a chosen robust approach in particular, affects the trial's final test and estimation decisions as well as the long-run performance of the design. The focus is on hypothesis testing procedures in the context of two-arm hybrid-control trial designs, and we analytically investigate and report on the connections between the Bayesian and the frequentist paradigm. Graphical and analytical tools are provided to compare test decisions induced by different approaches. Moreover, it is shown how test decisions can be adapted to satisfy frequentist constraints and/or data-adaptive costs of type I and type II errors under a fully Bayesian analysis.

C0381: Statistical inference for Levy-driven graph supOU processes: From short- to long-memory in high-dimensional time series

Presenter: **Almut Veraart**, Imperial College London, United Kingdom

Co-authors: Shreya Mehta

The aim is to introduce Levy-driven graph supOU processes. Such processes offer a parsimonious parametrization for high-dimensional time series, where dependencies between the individual components are governed by a graph structure. Specifically, a model specification is proposed that allows for a smooth transition between short- and long-memory settings while accommodating a wide range of marginal distributions. An inference procedure is further developed based on the generalized method of moments, its asymptotic properties are established, and its strong finite sample performance is demonstrated through a simulation study. Finally, the practical relevance of the new model and estimation method is illustrated in an empirical study of wind capacity factors in a European electricity network.

C0435: New allocation probability tests for improved power in response-adaptive clinical trials

Presenter: **Stina Zetterstrom**, University of Cambridge, United Kingdom

Co-authors: David Robertson, Sofia Villar

Response-adaptive clinical trials update the allocation probabilities to treatment arms based on the data collected so far in the trial. One objective for response-adaptive clinical trials is to ensure that patients in the trial will have a higher probability of getting the best treatment compared to using equal randomization. However, this imbalance in treatment allocation can lead to low power when testing for a treatment difference. Recent works propose a new testing approach that is based on the allocation probability (AP) instead of the outcome directly, which can increase the power. Alternative versions of the AP test are proposed, where the functional form of the test statistic is changed. The AP tests are evaluated in simulation studies, using the Bayesian response adaptive randomization (BRAR) algorithm for binary outcomes in a two-arm setting. The simulation studies show that the AP test can perform better in terms of power compared to traditional tests, while controlling type I error. Furthermore, changing the functional form of the AP test statistic can result in a higher power than the original AP test. While BRAR is used, the AP test can be used for any response-adaptive randomization, and its performance is studied in that intersection.

C0459: Delayed reveal randomization in platform trials: A simulation study

Presenter: **Rajenki Das**, MRC Biostatistics Unit, University of Cambridge, United Kingdom

Co-authors: Sofia Villar

Platform trials are increasingly adopted in intensive care and emergency medicine, comparing multiple interventions simultaneously. This presents a challenge of randomizing critically ill patients to numerous options at a sensitive point in their care. Consequently, "revealed randomization" (where treatment assignment is known immediately) is often preferred to reduce study burden. However, a delayed reveal, while practical for operational reasons, can complicate the rigorous delivery of randomization. A solution is presented through an example trial where participants are randomized to four treatment groups: Three active treatment + standard of care (SOC) arms, and one placebo + SOC arm. A unique aspect is that the fourth treatment group is reserved for a small fraction of the enrolled population meeting a specific eligibility criterion. Stratified randomization is employed, and the properties of this scheme are assessed, with particular emphasis on balance across treatment arms and strata. Findings illustrate a practical solution for a delayed-reveal randomization system that maintains robust performance while addressing the inherent challenges of multi-arm trials in acute care settings.

C0516: Why join professional organizations

Presenter: **Jessica Kohlschmidt**, The Ohio State University, United States

Life is in a busy, fast-paced world with many things on everyone's plate. Professional organizations are joined, so there is an opportunity to attend and present at conferences. It is also a line that can be added to a resume or CV. What other motivation is there for joining and volunteering with professional organizations, especially as young people? There are many organizations besides the large and well-known ones. There are many organizations available that are focused on an area or idea that may be of interest. More than one area of interest is often found that ignites interest and involvement with. There are sections of the ASA, regional chapters of the ASA, and independent organizations like the Caucus for Women in Statistics and Data Science. Personal experiences taking an active role in organizations and the value gained are shared.

CO112 Room BCB 403 CAUSAL INFERENCE FOR POLICY EVALUATION AND TREATMENT DECISION MAKING Chair: Nandita Mitra**C0861: How to ask the right question: Choosing estimands for policy research***Presenter:* **Nicholas Seewald**, University of Pennsylvania, United States

Evaluating the effects of policies is hard. Practical constraints like heterogeneous policy implementation details, slow rollouts, and small sample sizes add to the complexity of estimating causal effects. Adding to these difficulties is a lack of clearly articulated considerations for choosing estimands in policy evaluation research, which can lead to studies in which the choice of scientific question is guided by available analytic tools, rather than the other way around. The purpose is to discuss common estimands in policy evaluation, their strengths and limitations, and how they have become prominent in non-experimental research. More advanced causal questions are then highlighted, which better reflect the complex nature of health policy evaluation and require estimands beyond, say, the average treatment effect among the treated. The focus is then on complex health policy questions that are addressed with novel estimands and estimators to better understand policies nuanced effects and to move toward a healthier, more equitable future.

C0724: Causal transportability with local interference: A framework for policy effect prediction across regions*Presenter:* **Gary Hettinger**, New York University Grossman School of Medicine, United States*Co-authors:* Youjin Lee, Nandita Mitra

Policymakers often need to anticipate the impact of policies in regions where they were not yet implemented, relying on evaluations from regions where they have. For example, sweetened beverage taxes is studied in U.S. cities like Philadelphia, yet many municipalities are still considering adoption. Substantial heterogeneity in observed effects across regions underscores the need for methods that can reliably transport these effects to new settings. Transporting policy effects requires adjusting for demographic and contextual differences. However, existing approaches often overlook key factors like treatment interference from cross-border shopping and exposure heterogeneity from local economic dynamics. Transportability methods are extended by introducing a formal causal framework that incorporates spatial interference and heterogeneous exposure distributions, both of which may differ between observed and target regions. The method is assessed through simulation studies, and it is applied to estimate the impact of Philadelphia's beverage tax in other U.S. regions. Leveraging retail scanner data, zip-code-level Census demographics, and geographic proximity measures, tax effects are estimated in untreated areas, and the estimate's performance in already-taxed regions is assessed. Findings highlight the importance of accounting for regional variation not only in population characteristics, but also in spatial spillovers and economic dynamics, when transporting policy effects.

C0645: Bayesian evaluation of local policy effects on overdose outcomes*Presenter:* **Samrachana Adhikari**, NYU School of Medicine, United States

The United States drug overdose epidemic is a public health emergency resulting in a record number of overdose deaths. In response to the epidemic and to prevent and reduce harms from problematic drug use, many states have implemented overdose prevention and harm-reduction policies. In addition to the state laws, many cities have also created their own set of local policies specific to the conditions of the epidemic, unique to them. Despite the potential key role of local policies in reducing overdose outcomes, evidence of their population-level effectiveness is scant and contradictory, and there remain critical unanswered questions. A number of methodological challenges, due to small area data, co-occurring policies, heterogeneous policy environment, and exposure-confounder feedback, among others, have been major barriers to their evaluation. A Bayesian causal inference framework is presented to evaluate the impacts of local harm-reduction policies, adjusting for co-occurring policies and treatment-confounder feedback. The proposed longitudinal G-computation method with flexible Bayesian models incorporates random effects and adjusts for co-occurring policies through shrinkage priors to precisely estimate effects.

C0590: Novel Bayesian methods for sequential decision-making with incomplete information*Presenter:* **Arman Oganisian**, Brown University, United States

Chronic diseases are managed over time with a sequence of treatment decisions. Time-varying covariates are often used to tailor the decisions, but are usually monitored sporadically, leading to non-monotone missingness. In these settings, the interest is in estimating causal effects of different tailoring rules on a survival time outcome. It is formalized as a dynamic treatment regime (DTR) problem with a joint monitoring-treatment decision. Under specific identification assumptions, causal effects of such joint DTRs can be estimated, but require stratifying models across monitoring patterns. With many sparsely populated patterns, estimates can be unstable. Thus, Bayesian models are developed, equipped with a class of autoregressive priors that 1) smooth the models across time within a pattern and 2) smooth the models across monitoring patterns. The regularized models are embedded in a Bayesian g-computation procedure that draws from the posterior distribution of the causal effect of various DTRs. The model is applied to analyze survival among patients with pediatric acute myeloid leukemia (AML) enrolled in the AAML1031 trial. These patients move through a sequence of treatment courses in which a decision is made to withhold scheduled anthracycline chemotherapy (ACT). Since ACT is cardiotoxic, ejection fraction (EF) is sometimes - but not always - monitored to inform the withholding decision.

C0663: A flexible Bayesian g-formula for causal survival analyses with time-dependent confounding*Presenter:* **Liangyuan Hu**, Rutgers University, United States

In longitudinal observational studies with time-to-event outcomes, a common objective in causal analysis is to estimate the causal survival curve under hypothetical intervention scenarios. The g-formula is a useful tool for this analysis. To enhance the traditional parametric g-formula, an alternative g-formula estimator is developed, which incorporates the Bayesian additive regression trees (BART) into the modeling of the time-evolving generative components, aiming to mitigate the bias due to model misspecification. The focus is on binary time-varying treatments, and a general class of g-formulas is introduced for discrete survival data that can incorporate longitudinal balancing scores. The minimum sufficient formulation of these longitudinal balancing scores is linked to the nature of treatment strategies, i.e., static or dynamic. For each type of treatment strategy, posterior sampling algorithms are provided. Simulations are conducted to illustrate the empirical performance of the proposed method, and its practical utility is demonstrated using data from the Yale New Haven Health System's electronic health records.

CO056 Room BCB 405 BAYESIAN NONPARAMETRIC MIXTURE MODELING AND CLUSTERING (VIRTUAL) Chair: Catia Scricciolo**C1153: Enriched stochastic block models***Presenter:* **Louise Alamichel**, Bocconi University, Italy*Co-authors:* Daniele Durante, Tommaso Rigon, Francesco Gaffi

Stochastic block models learn group structures among nodes sharing similar connectivity patterns. Recent extensions have considered scenarios where multiple networks sharing the same nodes need to be analyzed jointly. State-of-the-art formulations typically assume that a unique partition of the nodes governs the clustering across all networks. In practice, the node partition describing block structures in one network may differ from that underlying another network, with the two often linked by refinements for specific purposes. For example, for two networks, the grouping of nodes in one layer may fragment into more detailed communities in another layer, reflecting different but related structures. To capture these architectures, we develop a novel enriched stochastic block model relying on two clustering structures, one for each network, and linked through a joint enriched Bayesian nonparametric prior. This construction induces a dependence across partitions, while retaining flexibility to model both shared and network-specific patterns. Inference under the proposed model is carried out via a collapsed Gibbs sampler on the cut posterior. Preliminary results on simulated data and in a study of summits co-attendances within a complex Mafia organization showcase the strengths of the proposed formulation, along with its ability to incorporate and learn relevant structures in networks.

C0283: Minimax-optimal dimension-reduced clustering for high-dimensional nonspherical mixtures*Presenter:* **Yuqi Gu**, Columbia University, United States

In mixture models, nonspherical (anisotropic) noise within each cluster is widely present in real-world data. Both the minimax rate and optimal statistical procedure for clustering under high-dimensional nonspherical mixture models are studied. In high-dimensional settings, the information-theoretic limits for clustering are first established under Gaussian mixtures. The minimax lower bound unveils an intriguing informational dimension-reduction phenomenon: There exists a substantial gap between the minimax rate and the oracle clustering risk, with the former determined solely by the projected centers and projected covariance matrices in a low-dimensional space. Motivated by the lower bound, a novel computationally efficient clustering method is proposed: Covariance projected spectral clustering (COPPO). Its key step is to project the high-dimensional data onto the low-dimensional space spanned by the cluster centers and then use the projected covariance matrices in this space to enhance clustering. Tight algorithmic upper bounds are established for COPPO, both for Gaussian noise with flexible covariance and general noise with local dependence. The theory indicates the minimax-optimality of COPPO in the Gaussian case and highlights its adaptivity to a broad spectrum of dependent noise. Extensive simulation studies under various noise structures and real data analysis demonstrate our method's superior performance.

C0901: On excess mass behavior in Gaussian mixture models with Orlicz-Wasserstein distances*Presenter:* **Nhat Pham Minh Ho**, University of Texas, Austin, United States

Dirichlet process mixture models (DPMM) in combination with Gaussian kernels have been an important modeling tool for numerous data domains arising from biological, physical, and social sciences. However, this versatility in applications does not extend to strong theoretical guarantees for the underlying parameter estimates, for which only a logarithmic rate is achieved. The aim is to (re)introduce and investigate a metric, named Orlicz-Wasserstein distance, in the study of the Bayesian contraction behavior for the parameters. It is shown that despite the overall slow convergence guarantees for all the parameters, posterior contraction for parameters happens at almost polynomial rates in outlier regions of the parameter space. The theoretical results provide new insight into understanding the convergence behavior of parameters arising from various settings of hierarchical Bayesian nonparametric models. In addition, an algorithm to compute the metric is provided by leveraging Sinkhorn divergences, and findings are validated through a simulation study.

C0571: Bayesian nonparametric mixture inconsistency for the number of components: How worried should we be in practice?*Presenter:* **Johan van der Molen Moris**, Pontificia Universidad Catolica de Chile, Chile*Co-authors:* Paul Kirk, Anthony Davison, Yannis Chaumeny

A Bayesian clustering approach is considered with the mixture of finite mixtures and Dirichlet process mixture models, popular due to uncertainty estimates for the number of clusters and efficient sampling methods. However, recent theoretical results show that Dirichlet process mixture models overestimate the number of clusters for large samples, and under misspecification, both models give inconsistent estimates. Furthermore, Bayesian mixture models give inconsistent numbers of clusters in some high-dimensional cases. In practice, Markov chain Monte Carlo summarization methods obtain a representative clustering for interpretation, and their effect on the number of clusters is not well studied. The consequences of summarisation methods are investigated for practical scenarios in light of asymptotic results with a simulation study and an application on gene expression data. Results show that, for the situations considered, the Dirichlet process mixture model leads to limited overestimation of the number of clusters for finite samples, which is corrected by some summarization methods. Misspecification leads to considerable overestimation of the number of clusters, but results are still interpretable. It is shown that certain summarization methods also lead to overestimation of the number of clusters, even for accurate estimates. For high-dimensional data, an illustration of the underestimation of the number of clusters is given, suggesting careful interpretation in practice.

C1042: A Bayesian method for learning mixture models of non-parametric components*Presenter:* **Yun Wei**, University of Texas at Dallas, United States*Co-authors:* Long Nguyen, Yilei Zhang, Aritra Guha

Mixture models are widely used in modeling heterogeneous subpopulations in data. Mixture models of parametric components (e.g., Gaussian mixture models) have been thoroughly studied on both statistical and algorithmic fronts. However, in the face of the increasing complexity of large-scale data, parametric assumptions such as Gaussianity are often unrealistic, and very little literature has been found on learning the mixture models of non-parametric components. In an effort to fill this gap, the identifiability issue in mixture models of non-parametric components is first addressed. Building on this, a framework is established using a mixture of Dirichlet processes to learn such models, and an efficient MCMC algorithm is developed to implement the method. The method can learn each component density without resorting to solving the mixing measure, thus providing a sample-efficient framework for learning subpopulation properties from data. The posterior contraction rate of the component density estimator of an almost polynomial order is also shown, which is a significant improvement from the logarithm convergence rate of solving mixing measures. This substantiates the sample efficiency and applicability of the method in learning non-parametric component densities.

CO020 Room BCB 406 ADVANCES IN MODELLING AND INFERENCE ON NETWORKS AND HYPERGRAPHS Chair: Nilanjan Chakraborty
C0773: Statistical inference for subgraph densities under induced random sampling from network data*Presenter:* **Ayoushman Bhattacharya**, Washington University in St. Louis, United States*Co-authors:* Nilanjan Chakraborty, Soumen Lahiri

Statistical inference is considered for subgraph densities of a general population network under without replacement sampling (SRSWOR). Under this sampling scheme, the asymptotic normality of the Horvitz-Thompson (HT) estimator is established for the population subgraph densities under minimal assumptions. The joint asymptotic normality of two subgraph densities is also established, which is crucial in establishing a weak convergence of the global transitivity of the sampled graph. To facilitate the inferential procedures, a jackknife and a bootstrap estimator of the unknown population variance are provided, and their consistency is established. Results find a useful application to the problem of testing the equality of two population graphs using the subgraph densities as the test statistic. Finally, a simulation study and a real data analysis are presented, which corroborate the theoretical findings.

C0959: Statistical inference for subgraph frequencies under exchangeable hyperedge models*Presenter:* **Nilanjan Chakraborty**, Missouri University of Science and Technology, United States*Co-authors:* Ayoushman Bhattacharya, Robert Lunde

The problem of statistical inference is considered for subgraph counts under an exchangeable hyperedge model. Various classes of subgraph statistics for hypergraphs are introduced, and inferential tools are developed for a notion of subgraph frequency that takes into account the multiplicity of an edge. It is further shown that a certain subclass of these subgraph statistics is robust to the deletion of low-degree nodes, facilitating inference in settings where low-degree nodes are more likely to be missing. A more traditional notion of subgraph frequency is also considered, which does not take into account multiplicity. It is demonstrated that while statistical inference based on limiting distributions for these statistics is possible in certain cases, a proper limiting distribution does not even exist in certain cases. The finite sample properties of the procedures are studied in both simulation studies and real-world datasets. For data analysis, new hypergraph datasets involving academic and movie collaborations are collected, and it is found that the inferential tools for hypergraphs have more power to distinguish between networks than traditional approaches based on binary adjacency matrices.

C0971: Network bootstrap using overlapping partitions*Presenter:* **Sayan Chakrabarty**, University of Michigan, United States*Co-authors:* Liza Levina

Bootstrapping network data efficiently is a challenging task. The existing methods tend to make strong assumptions on both the network structure and the statistics being bootstrapped, and are computationally costly. The aim is to introduce a general algorithm, SSBoot, for network bootstrap that partitions the network into multiple overlapping subnetworks and then aggregates results from bootstrapping these subnetworks to generate a bootstrap sample of the network statistic of interest. This approach tends to be much faster than competing methods, as most of the computations are done on smaller subnetworks. It is shown that SSBoot is consistent in distribution for a large class of network statistics under minimal assumptions on the network structure, and it is demonstrated with extensive numerical examples that the bootstrap confidence intervals produced by SSBoot attain good coverage without substantially increasing interval lengths in a fraction of the time needed for running competing methods.

C1079: Degree heterogeneity in higher-order networks: Inference in the hypergraph beta-model*Presenter:* **Sagnik Nandy**, University of Chicago, United States*Co-authors:* Bhaswar Bhattacharya

The aim is to introduce a multilayer hypergraph beta-model, extending the standard beta-model to networks with multi-way interactions and varying hyperedge sizes across layers. The statistical properties of this novel model are rigorously studied. First, the convergence rates of the maximum likelihood (ML) estimate are established, and their minimax optimality is proven. The ML estimate's limiting distribution is also determined, and asymptotically valid confidence intervals are constructed for model parameters. Next, the goodness-of-fit problem is addressed by analyzing the likelihood ratio (LR) test. Its asymptotic normality is derived under the null hypothesis, its detection threshold is determined, and its limiting power is analyzed. Notably, this detection threshold is minimax optimal, meaning no test can perform better below it. These theoretical findings are supported by numerical experiments. This not only develops a comprehensive framework for multilayer hypergraph beta-models but also addresses existing gaps in the graph beta-model literature, particularly concerning ML estimate minimax optimality and non-null LR test properties.

C1184: A technique for consistent response prediction on a collection of time series of graphs*Presenter:* **Aranyak Acharyya**, Johns Hopkins University, United States*Co-authors:* Francesco Sanna Passino, Michael Trosset, Carey Priebe

The aim is to propose a novel methodology for the response prediction task on a collection of time series of networks. The setting involves a collection of time series of networks, some of them labeled with responses, while most are unlabeled. Exploiting an underlying low-dimensional structure, the method predicts the response consistently at an unlabeled time series of interest. The applicability of the method is demonstrated on real data for studying biological circuits governing learning capability in larval *Drosophila*.

CO067 Room BCB 407 ADVANCES IN CAUSAL INFERENCE AND MEDIATION WITH COMPLEX DATA**Chair: Xizhen Cai****C0729: Causal mediation analysis for optimizing interventions using factorial designs***Presenter:* **Donna Coffman**, University Of South Carolina, United States

Optimization trials are used to design and build efficient, immediately scalable, and cost-effective interventions by selecting intervention components (i.e., any part of an intervention that may be separated out for experimental manipulation) for inclusion in an intervention. In most interventions, these components are designed to affect mediators (i.e., third variables) through which the components affect the outcome of interest. There is very little guidance on assessing mediation in the context of optimization trials. Methods are extended based on the potential outcomes framework for causal inference to draw more valid causal inferences about mediation in optimization trials that use a 2^K factorial design and effect coding. The focus is on defining, identifying, and interpreting causal mediation effects. For example, using effect coding, a mediated effect of a main effect may be interpreted as the effect of a factor on the outcome through the mediator, averaged over all combinations of levels of the remaining factors. These methods ultimately allow researchers to understand how their interventions work. With this knowledge, they are able to design more efficacious and cost-effective interventions by targeting the most relevant mediator(s). Understanding how interventions are effective is key to designing more powerful, efficient interventions.

C0467: Functional causal mediation analysis with zero-inflated count data*Presenter:* **Yeying Zhu**, University of Waterloo, Canada

Mediation analysis is crucial for understanding how a treatment exerts effects on an outcome via an intermediate variable, known as the mediator. Zero-inflated count outcomes and time-varying mediators are prevalent in fields such as biomedicine and biostatistics. To address the complex structure of the data, existing mediation analysis methodologies are extended by integrating a functional mediator in the context of zero-inflated count outcomes. The potential outcomes framework is employed to define the mediation effects of interest in this context and to provide the theoretical underpinning for our approach, including conditions for effect identification. Estimation and inference on the direct and indirect effects are performed by a quasi-Bayesian Monte Carlo approximation method using the well-known mediation formula. The methods are applied to study gender disparity in the number of readmission to ICU for patients in MIMIC-IV, an electronic health record database.

C0915: Comparisons of variable selection and inference methods in high-dimensional mediation analysis*Presenter:* **Xizhen Cai**, Williams College, United States*Co-authors:* Yuan Huang, Yeying Zhu

Mediation analysis is a framework for understanding how a treatment affects the outcome through intermediate variables, namely mediators. Over the past decades, large and high-dimensional datasets have become easily stored and publicly available. This has led to many recent advances in mediation analysis, including developing models to fit more complex data structures and methods for mediator selection in high-dimensional settings. The statistical inference procedure following the mediator selection is also an important step in the mediation analysis. The effect of different variable selection and inference procedures is studied through simulation studies. The simulation settings and the findings are discussed to provide guidelines that help distinguish among various approaches, highlight the advantages and disadvantages of each, and identify ones that perform better in certain scenarios.

C0396: Causal inference with outcome-dependent missingness and self-censoring*Presenter:* **Rohit Bhattacharya**, Williams College, United States*Co-authors:* Jacob Chen, Daniel Malinsky

Unobserved confounding and missing data are two of the most common issues that arise in observational studies. In particular, if the outcome variable directly affects its own missingness status, i.e., it is "self-censoring", the resulting non-ignorable missingness along with unobserved confounding may lead to severely biased causal effect estimates. A test is proposed based on a randomized incentive variable offered to encourage reporting of the outcome (e.g., randomized gift cards) that can be used to verify identification assumptions that are sufficient to correct for both self-censoring and confounding bias. Concretely, the test confirms whether a given set of pre-treatment covariates is sufficient to block all sources of confounding between the treatment and outcome, as well as all associations between the treatment and missingness indicator after conditioning on the outcome. It is shown that under these conditions, the causal effect is identified by using the treatment as a "shadow variable" that allows us to correct for self-censoring. This leads to an intuitive inverse probability weighting estimator that uses a product of the treatment and response weights. The efficacy of the test is evaluated, and the estimator is downstreamed via simulations.

C0975: Accounting for outcome spillover for causal inference with continuous spatiotemporal processes*Presenter:* **Duncan Clark**, Williams College, United States*Co-authors:* Conor Kresin, Martin Hazelton

The purpose is to demonstrate that, in a highly general parametric setting, causal inference with observational spatiotemporal data in the presence of arbitrary outcome spillover is feasible. A general framework is constructed for defining outcomes and associated causal estimands in continuous space-time, and it is demonstrated that estimation is achievable with a likelihood-based approach via stochastic expectation maximization. Leveraging results from point process theory, conditions necessary for estimation and subsequent inference are demonstrated. The proposed framework accommodates observational and experimental data, random and non-random treatment mechanisms, a general class of model specifications including those that allow for interaction between points, and general observation windows. The focus is pertinent to applications as diverse as epidemiology and finance, enabling previously impossible causal inference on rich continuous spatiotemporal data.

CO148 Room BCB 409 STATISTICAL METHODS FOR BIOMEDICAL APPLICATIONS**Chair: Satabdi Saha****C0222: Structured Bayesian variable selection for microbiome compositional data using graph-guided shrinkage***Presenter:* **Satabdi Saha**, University of Texas MD Anderson Cancer Center, United States*Co-authors:* Christine Peterson

A Bayesian regression model for microbiome compositional data that accounts for both sparsity and microbial structure is proposed. The method employs a structured horseshoe prior that encourages variable selection while borrowing strength across related taxa through a graph-informed shrinkage. To respect compositional constraints, regression coefficients are modeled under a sum-to-zero condition. Posterior inference is performed via an efficient blocked Gibbs sampler with Langevin updates. Through simulation studies and application to real microbiome data, improved feature selection and predictive performance are demonstrated over existing state-of-the-art algorithms.

C0228: DTH: A nonparametric test for homogeneity of multivariate dispersions*Presenter:* **Asmita Roy**, Johns Hopkins University, United States*Co-authors:* Glen Satten, Ni Zhao

Testing homogeneity across groups in multivariate data is an important scientific question in its own right, as well as an auxiliary step in verifying the assumptions of ANOVA. Existing methods either construct test statistics based on the distance of each observation from the group center or on the mean of pairwise dissimilarities among observations in a group. Both approaches can fail when the mean within-group distance is similar across groups, but the distribution of the within-group distances is different. This is a pertinent question in high-dimensional microbiome data, where outliers and overdispersion can distort the performance of a mean-dissimilarity-based test. The non-parametric distance-based test for homogeneity (DTH) is introduced, which measures dissimilarity between groups by comparing the empirical distribution of within-group dissimilarities using a combination of the Kolmogorov-Smirnov and Wasserstein distances. For more than two groups, pairwise group tests are combined using a permutation-based p-value. Through simulations, it is shown that the method has higher power than existing tests for homogeneity in certain situations and comparable power in other situations. A simple framework is also provided for extending the test to a continuous covariate.

C0766: Generalized odds: Distributional object to model physical activity*Presenter:* **Pratim Guha Niyogi**, University of Mississippi Medical Center, United States*Co-authors:* Kathryn Fitzgerald, Ellen Mowry, Vadim Zipunnikov

The advancement of mobile health (mHealth) and wearable technologies has provided valuable insights for medical and public health research, enabling a deeper understanding of the association between human behaviors and their impact on health. These data resources capture metrics such as physical activity (PA) and heart rate in near real-time, which have become increasingly popular due to their availability and wide range of applications in clinical research. Moreover, research on modeling the subject-specific distributional representations of such data has expanded considerably. The distribution representation of PA could be measured by the class of conditional and unconditional probability measures, such as survival and hazard functions, among many others. Generalized odds can provide a wide class of different distributional representations of the underlying data. This also produces a more nuanced interpretation of the relationships between PA and clinical outcomes. Motivated by the HEAL-MS study, a model that leverages scalar-on-distributional regression techniques is introduced to connect traditional measures of MS progression with this rich class of digital health data.

C0976: Towards scalable Gaussian process regression with a moment-based estimation procedure*Presenter:* **Savita Pareek**, Auburn University, United States*Co-authors:* Lionel Voirol, Roberto Molinari

Gaussian processes (GPs) are a powerful Bayesian non-parametric framework for supervised learning, enabling probabilistic function modeling and uncertainty quantification. By leveraging kernel functions, GPs effectively capture the underlying relationships in the data. However, likelihood-based estimation procedures for GPs incur a high computational cost of $O(n^3)$, making them impractical for large-scale datasets. To address this problem, various approximations have been proposed in the literature, such as global and local approximations (Nyström, subset selection, sparse kernels, etc.). Nevertheless, the estimation of these models is still limited by their significant computational cost in large-scale datasets containing millions or billions of data points. A novel moment-based estimation procedure is proposed based on the wavelet variance and the generalized method of wavelet moments (GMWM) estimator. This framework achieves a log-linear computational complexity of $O(n \log n)$, where n is the number of training points, making it highly scalable. The method is compared with the state-of-the-art implementations, such as Rs GauPro, and comparable predictive performance is observed with significant gains in computational time. Additionally, the method is applied to diverse real-world applications, and it is compared with existing software packages. The results demonstrate its practical scalability, establishing it as a promising solution for large-scale data modeling.

C0995: SMART-MC: Characterizing the dynamics of multiple sclerosis therapy transitions using a covariate-based Markov model*Presenter:* **Priyam Das**, Virginia Commonwealth University, United States

Treatment switching is common in managing multiple sclerosis (MS), where patients transition across disease-modifying therapies (DMTs) due to variable responses, disease progression, patient characteristics, or side effects. To study how covariates influence treatment transitions, a Markovian model sparse matrix estimation with covariate-based transitions in Markov chain modeling (SMART-MC) is adopted, in which transition probabilities depend on patient-level covariates. Modeling real-world transitions poses challenges, such as parameter identifiability and sparsity. Identifiability is addressed by constraining each transition-specific covariate vector to have unit L2 norm. Sparsely observed transitions are estimated as constants, and empirically unobserved transitions are set to zero, reducing complexity while preserving interpretability. To optimize the nonconvex likelihood, a scalable, parallel global optimization algorithm is developed, validated through benchmarks, and supported by theoretical guarantees. The analysis reveals distinct DMT switching patterns across MS subgroups defined by age, race, and clinical characteristics.

CC425 Room BCB 208 STATISTICAL COMPUTING AND SIMULATION**Chair: Bettina Gruen****C1081: Tidy simulation: Designing robust, reproducible, and scalable statistical simulations***Presenter:* **Erik-Jan van Kesteren**, Utrecht University, Netherlands

Monte Carlo simulation studies are at the core of the modern applied, computational, and theoretical statistical literature. Simulation is a broadly

applicable research tool, used to collect data on the relative performance of methods or data analysis approaches under a well-defined data-generating process. However, extant literature focuses largely on design aspects of simulation, rather than implementation strategies aligned with the current state of (statistical) programming languages, portable data formats, and multi-node cluster computing. Tidy simulation is proposed: A simple, language-agnostic, yet flexible functional framework for designing, writing, and running simulation studies. It has four components: A tidy simulation grid, a data generation function, an analysis function, and a results table. Using this structure, even the smallest simulations can be written in a consistent, modular way, yet they can be readily scaled to thousands of nodes in a computer cluster should the need arise. Tidy simulation also supports the iterative, sometimes exploratory nature of simulation-based experiments. By adopting the tidy simulation approach, researchers can implement their simulations in a robust, reproducible, and scalable way, which contributes to high-quality statistical science.

C1202: Automatic selection of hyper-parameters via the use of softened profile likelihood

Presenter: **Gengyang Chen**, University of Waterloo, Canada

Co-authors: Mu Zhu

The purpose is to extend a heuristic method for automatic dimensionality selection, which maximizes a profile likelihood to identify "elbows" in scree plots. The extension enables researchers to make automatic choices of multiple hyperparameters simultaneously. To facilitate the extension to multi-dimensions, a "softened" profile likelihood is proposed. Two distinct parameterizations of the solution are presented, and the approach is demonstrated on elastic nets, support vector machines, and neural networks.

C1326: Numerical generalized randomized Hamiltonian Monte Carlo for piecewise smooth target densities

Presenter: **Jimmy Huy Tran**, University of Stavanger, Norway

Co-authors: Tore Selland Kleppe

Traditional gradient-based sampling methods, like standard Hamiltonian Monte Carlo, require that the desired target distribution is continuous and differentiable. This limits the types of models one can define, although these models capture the reality in the observations better. In this project, generalized randomized Hamiltonian Monte Carlo (GRHMC) processes for sampling continuous densities with discontinuous gradient and piecewise smooth targets are proposed. The methods combine the advantages of Hamiltonian Monte Carlo methods with the nature of continuous time processes in the form of piecewise deterministic Markov processes to sample from such distributions. It is argued that the techniques lead to GRHMC processes that admit the desired target distribution as the invariant distribution in both scenarios. Simulation experiments verifying this fact and several relevant real-life models are presented, including a new parameterization of the spike and slab prior for regularized linear regression that returns sparse coefficient estimates and a regime-switching volatility model.

C1388: Scalable computations for generalized mixed effects models with crossed random effects using Krylov subspace methods

Presenter: **Fabio Sigrist**, ETH Zurich, Switzerland

Co-authors: Pascal Kuendig

Mixed effects models are widely used for modeling data with hierarchically grouped structures and high-cardinality categorical predictor variables. However, for high-dimensional crossed random effects, current standard computations relying on Cholesky decompositions can become prohibitively slow. The aim is to present novel Krylov subspace-based methods that address several existing computational bottlenecks. Among other things, various preconditioners are theoretically analyzed and empirically evaluated for the conjugate gradient and stochastic Lanczos quadrature methods, derive new convergence results, and computationally efficient methods are developed for calculating predictive variances. Extensive experiments using simulated and real-world data sets show that the proposed methods scale much better than Cholesky-based computations, for instance, achieving a runtime reduction of approximately two orders of magnitude for both estimation and prediction. Moreover, the software implementation is up to 10'000 times faster and more stable than state-of-the-art implementations such as lme4 and glmmTMB when using default settings.

C1254: Control variates for variance-reduced ratio of means estimators

Presenter: **Louison Bocquet-Nouaille**, ISAE-SUPAERO, France

Co-authors: Jerome Morio, Benjamin Bobbia

The control variates method is a classical variance reduction technique for Monte Carlo estimators that exploits correlated auxiliary variables without introducing bias. In many applications, the quantity of interest can be expressed as a ratio of expectations. The aim is to propose a new variance-reduced estimator for such ratios, which applies control variates to both the numerator and the denominator. The control variate coefficients are optimized jointly to minimize the variance of the resulting estimator. This approach guarantees variance reduction and naturally extends to approximate control variates. Simulation studies show significant variance reduction, particularly when the correlation between variables and control variates is strong. The practical value of the method is demonstrated through an aerospace case study.

CC439 Room BCB M201 ASSET PRICING AND TESTING

Chair: Robinson Kruse-Becher

C1249: Investor attention interaction and asset pricing: Evidence from co-search behavior

Presenter: **Jieying Zhang**, Tohoku University, Japan

Co-authors: Yasumasa Matsuda, Stanley Iat-Meng Ko, Runyu Dai, Chia Chun Lo

The aim is to investigate the explanatory and predictive power of investor attention interaction (IAI) in asset pricing, using 21.33 billion sequential searches from the U.S. Securities and Exchange Commission's EDGAR system. A time-varying, directional IAI network matrix is constructed based on investors' co-search behavior. Embedding this structure within a spatial arbitrage pricing theory (S-APT) framework, it is shown that the IAI network improves the performance of factor models by reducing pricing errors (alphas) to statistical insignificance. The inclusion of the IAI network reshapes risk exposures and return comovements across firms. It is further demonstrated that the IAI network predicts future firm fundamentals, including return on assets (ROA) and standardized unexpected earnings (SUE), and strong return predictability is exhibited. Extensive robustness checks confirm that these findings are not driven by alternative model assumptions, construction methods, or spurious linkages. Finally, evidence that the effect decays with network distance and that it decays even faster in degree-preserving placebo networks supports an attention-mediated information-diffusion mechanism. Results highlight the importance of incorporating investor attention dynamics into asset pricing models.

C1291: Exact single- and multi-currency tests of forward rate unbiasedness using Fama regressions

Presenter: **Richard Luger**, Laval University, Canada

Co-authors: Hsuan Fu

The aim is to develop exact simulation-based tests of the forward rate unbiasedness hypothesis (FRUH) at both the single- and multi-currency levels using the classical Fama regression. The procedures deliver finite-sample valid inference and are invariant to conditional heteroskedasticity, feedback from errors to future regressors, unequal sample lengths across currencies, and stretches of missing or structurally uninformative data. Within this framework, an exact global test is proposed together with single-step and step-down multiplicity adjustments that ensure strong family-wise error rate control when identifying currencies or parameters that deviate from the hypothesis. The methods are illustrated through a simulation study and an empirical application to an unbalanced panel of developed and emerging market currencies. Multiplicity-adjusted p-values reveal that rejections of the FRUH occur exclusively among emerging market currencies.

C0219: Testing model-based contributions in misspecified asset pricing models*Presenter:* **Cisil Sarisoy**, Federal Reserve Board, United States

Test statistics designed to evaluate the risk-premium contributions of asset pricing models are developed for individual assets. These statistics are robust to model misspecification, allowing for reliable use in both correctly and misspecified models. The method is applied to the well-known asset pricing models that include both traded and non-traded factors. The empirical results show that, in testing model-based contributions of asset pricing models, the misspecification-robust test statistics developed can lead to economically different conclusions from those derived under correctly specified models.

C0330: Pricing spread options with bivariate Gil-Pelaez inversion*Presenter:* **Daniel Miao**, National Taiwan University of Science and Technology, Taiwan

Pricing spread options remains challenging due to their dependency on the difference between two underlying asset prices. Existing numerical approaches, notably the two-dimensional FFT method introduced by another study, are widely applied but sensitive to grid and damping parameter choices. To address these limitations, a pricing method based on the bivariate Gil-Pelaez inversion formula is proposed, leveraging the joint characteristic function of the asset distributions to directly evaluate a two-dimensional integral without extensive parameter tuning. Numerical experiments across diverse models, including geometric Brownian motion (GBM), jump-diffusion (JD), variance gamma (VG), stochastic volatility (SV), and stochastic volatility with jumps (SVJ), demonstrate the method's superior accuracy and numerical stability compared to FFT-based approaches. Additionally, cubic spline interpolation is integrated to significantly reduce computational complexity. The interpolation approach achieves equivalent accuracy to dense numerical grids but requires substantially fewer integration points, thereby notably shortening computation time. Results underline the effectiveness and computational advantages of the bivariate Gil-Pelaez inversion method, particularly for complex financial models involving heavy tails and jumps.

CC404 Room BCB 408 APPLIED STATISTICS**Chair: Silvia Montagna****C1487: Geographically weighted regression for air quality low-cost sensor calibration***Presenter:* **Jean-Michel Poggi**, University Paris-Saclay Orsay, France*Co-authors:* Bruno Portier, Emma Thulliez

The focus is on the use of Geographically Weighted Regression (GWR) to correct low-cost air quality sensor measurements. Those sensors are of major interest in the current era of high-resolution air quality monitoring at the urban scale, but require calibration using reference analyzers. The results for NO₂ are provided along with comments on the estimated GWR model and the spatial content of the estimated coefficients. The study has been carried out using the publicly available SenseURCity dataset in Antwerp, which is particularly relevant because it includes 9 reference stations and 34 micro-sensors, all collocated and deployed within the city.

C1201: Tree-guided variable selection methods for the regularized Dirichlet-tree multinomial regression*Presenter:* **Alysha Cooper**, University of Guelph, Canada*Co-authors:* Zeny Feng, Ayesha Ali

Hispanics and Latino Americans make up one of the largest minority groups in the United States, yet it is underexplored why they face a disproportionately high prevalence of chronic diseases such as asthma, diabetes, and obesity. While socioeconomic disadvantages and limited access to healthcare are known contributors, less attention has been given to how lifestyle changes following relocation to the U.S. may influence gut health and, in turn, overall health. Gut microbiome composition can be measured as bacterial counts organized along a taxonomic tree. The Dirichlet-tree multinomial (DTM) regression model accommodates the high variability in microbial count data while respecting evolutionary relationships among taxa. Variable selection becomes critical in DTM regression due to the high dimensionality of the outcome and many potentially relevant factors influencing gut bacteria. A key issue is that standard selection approaches may ignore dependencies among outcomes arising from the tree structure. Two tree-guided penalties for DTM regression are proposed: 1) the tree-guided sparse group lasso and 2) the tree-guided hierarchical lasso. Both penalties introduce sparsity into the model while leveraging known relationships among taxa in the tree. The regularized DTM regression model is demonstrated through simulations and analysis of gut microbiome data from the Hispanic Community Health Study/Study of Latinos.

C0206: Beyond the sample: Bootstrap ex ante accuracy estimation under any regression model*Presenter:* **Alicja Wolny-Dominiak**, University of Economics in Katowice, Poland*Co-authors:* Tomasz Zadło

Bootstrap methods are widely used in data analysis for various purposes, including the estimation of sampling distributions, assessing the precision of the estimation, estimating and correcting bias, and improving hypothesis testing. They are also popular in time series and spatial prediction, utilizing both cross-sectional and longitudinal data, to estimate the ex-ante prediction accuracy. However, the model-based bootstrap algorithms currently available are designed for specific parametric models. We propose a new bootstrap algorithm that can be applied to any model, enabling the generation of both sample and out-of-sample data. This allows for the estimation of ex-ante prediction accuracy for forecasts derived from both parametric and nonparametric models. Furthermore, our approach facilitates the estimation of ex-ante prediction accuracy under machine learning models, where traditional methods, such as k-fold cross-validation, assess the ex-post prediction accuracy. Our proposal is extensively tested in two Monte Carlo simulation studies. In the first study, we analyse its performance under a linear mixed model, comparing it with traditional bootstrap methods, such as parametric and residual bootstrapping, which have been shown to possess good properties in this case. The second simulation study examines the properties of our method under a nonparametric model, where estimating ex-ante prediction accuracy presents a significant challenge.

C1499: Tree-based learning of structural changes in mixture models with uncertainty for rating data*Presenter:* **Carmela Cappelli**, University of Naples Federico II, Italy*Co-authors:* Rosaria Simone

Qualitative data on individuals' perceptions and subjective evaluations are collected through official surveys using rating scales, often administered repeatedly over multiple time waves. In many cases, only aggregated response distributions are made publicly available, with no access to individual-level data. In this context, we combine Atheoretical Regression Trees with CUB models to analyze rating data and identify structural changes over time in the main characteristics of the response distributions. The chosen modeling framework is well-suited to parameterize both the latent feeling and the uncertainty underlying the ordinal evaluation, with a time-varying structure. Then, the use of Atheoretical Regression Trees involves introducing an artificial covariate that preserves temporal ordering, enabling the segmentation of feeling and uncertainty measures into homogeneous time intervals, which can then be interpreted in light of political and socio-economic events. We illustrate the proposed approach using data from the Consumer Confidence Survey issued monthly by the Italian National Institute of Statistics (ISTAT), assessing whether and to what extent individuals changed their judgments and expectations regarding price levels from 1994 to 2019. The proposed method demonstrates satisfactory performance when compared with the Bai and Perrons test for structural change detection, which is used as the benchmark.

Sunday 14.12.2025

08:45 - 10:25

Parallel Session F – CFE-CMStatistics 2025

CI011 Room BCB 307 ADVANCES IN BAYESIAN COMPUTATION FOR COMPLEX MODELS**Chair: Matias Quiroz****C0170: Bayesian nonparametric spectral analysis of locally stationary time series****Presenter:** Renate Meyer, University of Auckland, New Zealand**Co-authors:** Yifu Tang, Claudia Kirch, Kate Lee

Based on a novel dynamic Whittle likelihood approximation for locally stationary processes, a Bayesian nonparametric approach to estimating the time-varying spectral density is proposed. This dynamic frequency-domain-based likelihood approximation characterizes the time-frequency evolution of the process by utilizing moving periodograms previously introduced in the bootstrap literature. The posterior distribution is obtained by updating a bivariate extension of the Bernstein-Dirichlet process prior with the dynamic Whittle likelihood. Asymptotic properties such as sup-norm posterior consistency and L2-norm posterior contraction rates are presented. Additionally, this methodology enables model selection between stationarity and non-stationarity based on the Bayes factor. The finite-sample performance of the method is investigated in simulation studies, and applications to real-life data sets are presented.

C0171: A modified algorithm for MCMC in large dimensional models**Presenter:** Michael Pitt, Kings College London, United Kingdom

Large-dimensional models are of interest in classical and Bayesian inference. Bayesian approaches, based upon Markov chain Monte Carlo (MCMC), can run into difficulties when the parameter dimension becomes large. The purpose is to introduce a simple but effective remedy by modifying existing algorithms. This allows MCMC to work effectively in large dimensions. Asymptotic results are provided, which simplify the calculation and optimization of the integrated autocorrelation time (IACT). These asymptotic results and associated central limit theorems (CLTs) are illustrated throughout with examples. It can be seen that the results are robust in that they hold even when the number of parameters is modest (at say 30).

C0172: Bayesian inference for nonlinear and non-Gaussian state-space models with global-local shrinkage process priors**Presenter:** Mattias Villani, Stockholm University, Sweden

The recently proposed global-local shrinkage process priors for time-varying parameter models allow parameters to be essentially constant for longer spells, followed by periods of rapid change or jumps. Most of this literature uses linear and conditionally Gaussian models. The standard posterior samplers based on particle methods, typically used for nonlinear and non-Gaussian state space models, struggle with the parameter constancy implied by the global-local shrinkage processes. Several alternative algorithms are presented, and their performance is illustrated on some commonly used time-varying parameter models in statistics and econometrics.

CO052 Room BCB G07 HiTEC: ADVANCES IN ECONOMETRIC METHODS FOR ECONOMICS AND FINANCE Chair: Alessandra Amendola**C0658: Pink queue and double standard: A new assessment of the gender gap in the Italian academia****Presenter:** Marianna Brunetti, University of Rome Tor Vergata, Italy**Co-authors:** Annalisa Fabretti, Mariangela Zoli

The purpose is to study the gap in career progression in Italian academia based on a discrete-time Markov chain, which has been largely used in demography and the labor market, but never to assess potential gaps in career progressions in the university system. The probability of achieving different academic positions is estimated starting from entry-level roles, and the average time required for career advancement (i.e., promotion to a higher role), and to assess potential gaps in these metrics, the chain is applied on the subset of male and female professors, controlling for influential factors such as productivity, generation, and academic age, where several findings are reported. First, the gap in the probability of career advancement is marked over a 10-year horizon and is only partially absorbed over a 40-year full-career horizon. Second, female scholars wait up to 3 years longer than their male scientific counterparts. This phenomenon is labeled as "pink queue". Last, scholars are classified based on the quartile of the distribution they belong to in terms of h-index, and it is found that the likelihood of becoming a full professor for females in the top (third) quartile of the h-index distribution is in line with the same chances estimated for males in the third (second) quartile of the distribution. This result is labeled as "double standard" as it seems to suggest that women must achieve higher standards with respect to males in order to get to the same career level.

C0392: Systemic risk in the European insurance sector**Presenter:** Giovanni Bonaccolto, University of Enna Kore, Italy**Co-authors:** Nicola Borri, Andrea Consiglio, Giorgio Di Giorgio

The dynamic interdependencies are investigated between the European insurance sector and key financial markets-equity, bond, and banking-by extending the generalized forecast error variance decomposition framework to a broad set of performance and risk indicators. The empirical analysis, based on a comprehensive dataset spanning January 2000 to October 2024, shows that the insurance market is not a passive receiver of external shocks but an active contributor in the propagation of systemic risk, particularly during periods of financial stress such as the subprime crisis, the European sovereign debt crisis, and the COVID-19 pandemic. Significant heterogeneity is observed across subsectors, with diversified multiline insurers and reinsurance playing key roles in shock transmission. Moreover, the granular company-level analysis reveals clusters of systemically central insurance companies, underscoring the presence of a core group that consistently exhibits high interconnectivity and influence in risk propagation.

C1190: Long-term forecasting of stock returns: Avoid overly complex machine learning and prioritize benchmarking**Presenter:** Parastoo Mousavi, Bayes business school, United Kingdom**Co-authors:** Jens Perch Nielsen, Tatiana Franus

The aim is to investigate long-term stock return forecasting, emphasizing the importance of systematic benchmarking of both dependent and independent variables. It is shown that adjusting independent variables for relevant benchmarks substantially improves predictive performance, a largely unexplored approach in the literature. Using a range of machine learning methods, it is further demonstrated that simple, carefully designed models provide more reliable forecasts than automated approaches in data-constrained environments.

C0687: On the importance of ESG pillars in asset allocation: Evidence from simulated portfolios**Presenter:** Alessandra Amendola, University of Salerno, Italy**Co-authors:** Luigi Aldieri, Vincenzo Candila

Over the past two decades, corporate Environmental, Social, and Governance (ESG) pillars have received increasing attention from both theoretical and institutional perspectives. In principle, the asset allocation process is expected to favor firms with stronger ESG profiles, as both private and institutional investors are assumed to incorporate ESG factors into their investment decisions. However, empirical evidence on the actual impact of ESG considerations on portfolio choices remains mixed, and the effectiveness of ESG integration continues to be debated in the literature. The contribution to the debate is by analyzing the composition of optimal portfolios, each consisting of 30 assets equally distributed across ESG rating categories: 10 Low, 10 Medium, and 10 High-rated stocks. To ensure robustness, 10,000 such portfolios are simulated, each built from a random selection of the top 200 NASDAQ constituents by market weight. Conditional covariance matrices are estimated using the dynamic conditional

correlation (DCC) model, and portfolio weights are derived according to the global minimum variance, maximum diversification, and maximum decorrelation criteria, along with the benchmark equally weighted portfolio.

CO107 Room BCB G08 APPLICATIONS IN HIGH-DIMENSIONAL DATASETS
Chair: Alessia Paccagnini
C1229: On the identification of long and short memory components with applications to volatility

Presenter: **Tommaso Proietti**, University of Roma Tor Vergata, Italy

Co-authors: Alessandra Luati, Shelton Peiris, Gnanadarsha Dissanayake

The purpose is to develop a methodology for deriving the autoregressive and moving average coefficients of long memory models, with a focus on ARFIMA and Gegenbauer processes. The approach is based on the cepstrum of the long memory process, which allows for a systematic decomposition of the impulse response function into long memory and short memory components. This decomposition provides a separation between persistent dynamics and short-run fluctuations. The method is compared to existing approaches, such as the Beveridge-Nelson decomposition, and its application is illustrated in analyzing volatility processes and their roughness.

C1275: A new model to forecast energy inflation in the euro area

Presenter: **Mario Porqueddu**, European Central Bank, Germany

Co-authors: Marta Banbura, Elena Bobeica, Alessandro Giammaria, Josha van Spronsen

Energy inflation is a major source of headline inflation volatility and forecast errors; therefore, it is critical to model it accurately. The aim is to introduce a novel suite of Bayesian VAR models for euro area HICP energy inflation, which adopts a granular, bottom-up approach disaggregating energy into subcomponents, such as fuels, gas, and electricity. The suite incorporates key features for energy prices: Stochastic volatility, outlier correction, high-frequency indicators, and pre-tax price modelling. These characteristics enhance both in-sample explanatory power and forecast accuracy. Compared to standard benchmarks and official projections, the BVARs achieve better forecasting performance, particularly beyond the very short term. The suite also captures a sizable variation in the impact of commodity price shocks, pointing to higher elasticities at higher levels of commodity prices. Beyond forecasting, the framework is also useful for scenario and sensitivity analysis as an effective tool to gauge risks.

C1368: Multiresolution analysis of climate risk spillovers

Presenter: **Fabio Parla**, University of Palermo, Italy

Co-authors: Andrea Cipollini, Iolanda Lo Cascio

The focus of this study is on the climate risk implications for financial stability using the Diebold-Yilmaz connectedness approach. In particular, a methodology of another study is employed based on the computation of the marginal contribution to R². For this purpose, data for Europe and the US on Exchange-Traded Funds (ETFs) tracking "green" and "brown" stock portfolios are used. ETFs are an innovative financial instrument providing benefits to investors in terms of diversification (by tracking a basket of stocks), and given their broad coverage, they also typically feature low management fees and transaction costs. Contribution to the literature is twofold. First, while previous studies on ETF spillovers have focused only on world stock portfolios, the roles played by Europe and the US (tracking different environmental regulation policies) are disentangled. Moreover, using the multiresolution analysis approach put forward by the prior study, the focus is on the analysis of (total and directional) connectedness in a time-frequency framework. More specifically, the analysis studies ETF spillovers across frequency bands (associated with different investment horizons) and also examines their evolution over time.

C1392: A network analysis of success at the Eurovision song contest

Presenter: **Alessia Paccagnini**, University College Dublin, Ireland

Co-authors: Alessia Morrone, Barbara Bedowska-Sojka, Sabrina Giordano, Claudia Tarantola

What determines the success of a song at the Eurovision Song Contest? Is it driven by individual taste and geopolitical alliances, or do shared musical structures shape audience preferences? The purpose is to use the Eurovision Song Contest as a natural laboratory to explore the complex interdependencies between music, culture, and politics through the lens of network analysis. Using a purpose-built dataset, two types of networks are constructed: undirected similarity networks connecting songs based on musical, lyrical, and performative features, and bipartite networks linking voting countries to the songs they voted for. Centrality measures and clustering algorithms are applied to uncover structural patterns within these networks. Findings reveal that winning entries are often not the most central in the similarity networks. Instead, they tend to occupy peripheral positions, distinguishing themselves through unconventional features. At the same time, recurring characteristics emerge among highly voted songs: the use of English lyrics, minor keys, medium-to-high tempo (BPM), and strong visual performances. Furthermore, unexpected voting affinities are detected between countries with limited historical or cultural ties. It is illustrated how network-based approaches can uncover latent structures in cultural phenomena, offering an original perspective on collective behavior, aesthetic preferences, and strategic alliances in an international contest.

CO073 Room BCB G09 MACROECONOMIC POLICIES AND MACROECONOMETRICS
Chair: Etsuro Shioji
C0532: Tweets on childbirth in Japan: Sentiment and influence

Presenter: **Toshitaka Sekine**, Hitotsubashi University, Japan

Against the backdrop of Japan's dangerously low fertility rate, sentiment analysis is conducted on Japanese tweets about childbirth using the recent advances of NLP. It is found that these posts are associated with positive arousal (mentally activated) and negative valence (unpleasant). It is also found that more significantly negative posts (higher arousal and lower valence) are more influential.

C0539: Assessing partial exogeneity of shock proxies: A Bayesian proxy-DSGE framework

Presenter: **Ryohei Oishi**, University College London, United Kingdom

The proxy DSGE framework is introduced, a Bayesian econometric methodology for evaluating the validity of externally identified structural shock proxies using dynamic stochastic general equilibrium (DSGE) models. The framework assesses proxies based on relevance and partial exogeneity conditions, which is a weaker but more empirically tractable notion of exogeneity that requires proxies to be uncorrelated with non-targeted structural shocks within a DSGE model. These conditions are econometrically implemented through a spike-and-slab prior, which allows posterior loadings on shocks to be exactly zero with positive probability. The framework further provides a coherent estimation strategy for DSGE models with potentially weak proxies by extracting informative signals from their joint distribution. Applying the framework to a medium-scale DSGE model with US monetary policy proxies shows that the Romer-Romer narrative proxy and several high-frequency identified proxies are valid, although the latter are weak. Incorporating these proxies yields a more inertial Taylor rule and more precise estimates of monetary policy shocks.

C0470: Historical debt dynamics in the U.S.: Evidence from a TVPVARX model

Presenter: **Hiroshi Morita**, Institute of Science Tokyo, Japan

The drivers of fluctuations in U.S. government debt are examined using long-run data from 1880 to 2023. It is empirically assessed whether debt dynamics can be attributed to a systematic fiscal rule, fiscal policy shocks, and other identified structural shocks within a time-varying parameter VAR framework incorporating a debt-feedback rule of a prior study. The model treats the real GDP growth rate, the real long-run interest rate, and the primary surplus-to-GDP ratio as endogenous variables, with the debt-to-GDP ratio specified as exogenous. Debt dynamics are then derived explicitly from these three variables via the government budget constraint. Empirical findings are as follows: (1) discretionary fiscal policy

shocks were the main contributors to debt accumulation during the 1940s through the 1960s, whereas business cycle shocks have dominated since the Global Financial Crisis; (2) the fiscal stance measured by the elasticity of the primary surplus with respect to debt is positive throughout the sample, indicating Ricardian behavior by the U.S. government (3) although debt responses to shocks have evolved smoothly over time, the difference between the responses at the beginning and end of the sample period is statistically significant. These results shed new light on the historical drivers of U.S. debt dynamics and offer a framework for evaluating future fiscal sustainability.

C0530: Is the market for the Japanese government bonds still insulated from external forces?

Presenter: **Etsuro Shioji**, Chuo University, Japan

The purpose is to examine whether the market for Japanese government bonds (JGBs) is becoming integrated with the global financial market. Historically, the market has been dominated by domestic players, who often hold on to the bonds due to institutional or regulatory considerations. This tendency has effectively insulated it from external pressures. However, with the increasing presence of foreign investors, who may not hesitate to dump the bonds as they see fit, the market may be becoming more susceptible to events in external financial markets. Investigating this issue is complicated by the yield curve control (YCC) policy implemented by the Bank of Japan (BOJ) between 2016 and 2024. Under this policy, the BOJ actively intervened to stabilize yields, thereby suppressing much of their reaction to outside shocks. The aim is to address this challenge by employing an improved version of the bond market pressure index, which was developed in another study. This approach is based on the premise that, under the YCC regime, market pressures would manifest themselves in the amounts of intervention conducted by the BOJ. It is estimated how the index has responded to external shocks such as US monetary policy, and it is examined whether the reaction has strengthened over time.

CO029 Room Virtual R01 BAYESIAN NONPARAMETRICS FOR TIME SERIES ANALYSIS

Chair: Ramses Mena

C0751: Cluster analysis for longitudinal data

Presenter: **Asael Fabian Martínez Martínez**, Universidad Autonoma Metropolitana, Mexico

Co-authors: Ivonne Ramirez-Silva, Ruth Fuentes-Garcia

The identification of latent profile trajectories in longitudinal studies represents an important challenge for specialists. Many of the statistical methodologies are based on growth curve or mixed-effects models, and, for Bayesian nonparametric methods, Dirichlet process mixture models are widely used together. A clustering methodology is presented for longitudinal data based on mixture models generated by a discrete random probability measure whose weights are decreasingly ordered by construction. A straightforward procedure to merge some estimated groups is also provided, since it could happen that there are many of them to be easily explained by experts. The methodology is illustrated using simulated and real datasets.

C0944: Beta-CoRM binary matrix factorization

Presenter: **Jose Perusquia**, Universidad Nacional Autonoma de Mexico, Mexico

Co-authors: Jim Griffin, Cristiano Villa

Binary data arises naturally across diverse fields, including psychology, natural language processing, and computer science. These datasets are often high-dimensional, with many features per observation, motivating the need for compact, low-dimensional representations. Traditional variable selection techniques may discard important information, whereas matrix factorization methods provide a more flexible alternative by uncovering latent structures. The aim is to propose a Bayesian nonparametric binary matrix factorization model tailored to grouped binary data, a setting often overlooked in the literature. The model assumes that each binary feature is driven by a finite set of latent traits. Two constraints are imposed: One matrix encodes the presence or absence of latent traits as binary values, while the other consists of non-negative entries constrained to the unit interval. A key innovation is the use of a link function that introduces a second layer of binary latent indicators, modulating the generation of observed binary features. This hierarchical structure allows for interpretable and flexible modelling of complex binary datasets.

C0973: Distribution splicing via Bayesian non-parametric methods

Presenter: **Jorge Gonzalez Cazares**, IIMAS-UNAM, Mexico

Co-authors: Martin Bladt

An objective-based Bayesian non-parametric framework for studying time-to-event data is proposed, where the prior distribution is allowed to depend on an additional random source, and may update with the sample size. Such scenarios are natural, for instance, when considering empirical Bayes techniques or dynamic expert information. Conditionally inhomogeneous independent increment processes are used with non-decreasing sample paths. The asymptotic behavior is studied by showing that Bayesian consistency and Bernstein-von-Mises theorems may be recovered under suitable conditions on the asymptotic negligibility of the stochastic prior sequences. The non-asymptotic behavior of the posterior is also considered. Namely, upon conditioning, an efficient and exact simulation algorithm for the paths of the beta Levy process is provided. As a natural application, it is shown how the model can provide an appropriate definition of non-parametric spliced models, targeting data where an accurate global description of both the body and tail of the distribution is desirable. The Bayesian non-parametric nature of the proposed estimators can offer conceptual and numerical alternatives to their parametric counterparts.

C1198: Bayesian nonparametric mixtures of Markov processes

Presenter: **Ramses Mena**, Universidad Nacional Autonoma De Mexico, Mexico

The purpose is to introduce a framework for constructing nonparametric Markov processes by starting from a base process with well-defined transition probabilities and invariant distributions. A random mixture mechanism is then used to generate new processes, where the dynamics are governed by random coefficients derived from the mixture structure. This approach yields a flexible class of Markov models that preserve key invariance properties while allowing for greater adaptability in modeling complex systems. To fully develop the framework, it is necessary to study fundamental theoretical aspects such as stability and symmetry. In addition, because the transition mechanism might involve infinitely many components, new strategies for estimation and prediction are required.

CO077 Room BCB 206 ARDL MODELS AND COINTEGRATION: A TRIBUTE TO H. PESARAN

Chair: Kleanthis Natsopoulos

C0667: A spatiotemporal autoregressive factor model of the global business cycle

Presenter: **Matthew Greenwood-Nimmo**, University of Melbourne, Australia

Co-authors: Tomohiro Ando, Matthew Greenwood-Nimmo, Yongcheol Shin, Chaowen Zheng

To study the synchronicity of business cycles across countries, a new heterogeneous parameter panel data model is developed, in which the global business cycle is characterized as a spatiotemporal autoregressive process with a common factor error structure. A modified quasi-maximum likelihood approach is developed to estimate the model parameters in the presence of parameter heterogeneity and endogeneity. It is proven that the estimators are consistent and asymptotically normally distributed, and Monte Carlo simulations are used to show that their finite-sample performance is satisfactory. A framework is then developed for network analysis based on forecast error variance decomposition and diffusion multipliers. These tools are applied to analyze business cycle synchronization among 79 countries over the 50-year period 1970-2019.

C1139: Hawks, doves and asymmetry in US monetary policy: Evidence from a dynamic quantile regression model

Presenter: **Yongcheol Shin**, University of York, United Kingdom

The aim is to develop a quantile error correction representation of US monetary policymaking over the period 1964q to 22007q2. The model is unique as it is capable of simultaneously modelling three distinct forms of asymmetry: long-run (reaction) asymmetry, short-run (adjustment)

asymmetry, and quantile-specific (locational) asymmetry. Results reveal that: (i) the reaction function is linear in both output and inflation gaps in the lower conditional quantiles of the interest rate and that the Taylor principle is not upheld; (ii) in the mid quantiles, the reaction function exhibits a good performance bias whereby the Taylor principle is observed for positive inflationary shocks only; and (iii) no stable long-run relationship is observed in the upper quantiles. This final result is interpreted in terms of a modified Goodhart's Law, whereby excessive policy aggression undermines observed empirical regularities.

C0855: Advancing ARDL capabilities in R: Methods, extensions, and usability enhancements

Presenter: **Kleanthis Natsiopoulou**, University of Dundee, United Kingdom

The ARDL package in R, widely adopted by researchers, institutions, and central banks, has played a key role in facilitating applied cointegration analysis in an open-source setting. Building on this foundation, future developments aim to expand the package's methodological scope and practical usability. Planned extensions include the implementation of several advanced model variants within the ARDL framework, along with improvements to the structure and performance of the codebase. Efforts are also underway to explore the development of an interactive user interface to make the functionality more accessible to a broader audience. By advancing methodological coverage and improving usability, the goal is to support both research and real-world applications, while encouraging replicable results in time series econometrics.

C0794: ARDL models with high frequency intervals and dynamic sampling: A case study from crypto liquidity

Presenter: **Daniel Finnan**, Conservatoire national des arts et metiers (CNAM), France

Cointegration analysis and ARDL models are typically applied to quarterly or monthly data, with inherently less noise and variance than more granular intervals. However, the bounds testing procedure at the heart of ARDL analysis can equally be applied to higher frequency data if the model being explored has a strong theoretical grounding. Furthermore, with a greater sample size, techniques employing dynamic sub-samples can be employed to capture changes in the cointegrated relationship, using rolling windows and recursive sampling. This is akin to testing for structural breaks, but places the emphasis on the F- and t-statistics in the bounds testing procedure, rather than traditional tests for stability such as CUSUM. It also has the advantage of enabling the econometrician to build up a picture of changes in the speed of adjustment and long-run relationship over time. When applied to time series data using a variety of frequencies, it is possible to establish results giving both an overall view and a much more detailed, fine-grained perspective of changes to a relationship. An illustration of the approach is presented in the form of liquidity flows on cryptocurrency exchanges, using both daily and hourly intervals, over a year, month-by-month, and using dynamic sampling.

CO172 Room BCB 209 STATISTICAL MODELS FOR SHADOW ECONOMY AND FRAUD DETECTION

Chair: Maria Felice Arezzo

C0272: Fragility and resilience: Industrial structure under the threat of dimensional traps

Presenter: **Raffaella Coppier**, University of Macerata, Italy

Size-dependent fiscal policies are suspected to lead to sub-optimal outcomes in terms of fighting evasion and efficient allocation of investments by firms. A prior study describes how size-dependent tax enforcement may be a source of a dimensional trap in a single-firm evasion model. Their results suggest that, depending on the policy parameters, it may be optimal for the firm not to invest and remain small to face a lower probability of being monitored and punished in case of evasion. The aim is to move to a heterogeneous setting, thus considering industrial structures composed of many firms with different dimensions. When many firms are considered, the consequences of size-dependent policies and potential dimensional traps on the industrial structure and its evolution over time must be studied. A non-deterministic discrete-time dynamic model is proposed. By combining analytical findings and simulation results, it is shown that an unwise choice of policy parameters may determine, in the long run, an industrial structure characterized by a small number of large firms and a plethora of small firms, with the latter being marked by inefficient resource allocation and non-compliant behavior with regard to tax regulations. These results are robust across different choices of the initial industrial structure.

C0319: The potential of Benford's law for detecting tax fraud at the firm level in Italy

Presenter: **Luisa Scaccia**, University of Macerata, Italy

Co-authors: Raffaella Coppier, Elisabetta Michetti

The purpose is to explore the potential of Benford's Law (BL) to detect corporate tax evasion, focusing on Italy, where tax non-compliance remains a major economic issue. BL describes the expected frequency of leading digits in naturally occurring numerical data, and significant deviations may signal irregularities such as tax evasion. BL is applied to financial statement data from the AIDA database (2018-2022). A firm is considered potentially non-compliant if its data deviates from BL in at least one year. This offers a firm-level indicator that can guide audits and act as a proxy for tax evasion in empirical research. Such a proxy helps address the lack of firm-level data, enabling studies on how tax evasion affects firm growth or market distortions. To assess its validity, the distribution of non-compliant firms (as flagged by BL) is compared across regions, sectors, and firm sizes with official 2022 data from the Italian Ministry of Economy and Finance and the Revenue Agency, based on Synthetic Tax Reliability Indices (ISA). Results show a strong alignment between BL-based detection and official indicators, highlighting the usefulness of BL as a low-cost, data-driven tool for identifying and studying tax evasion at the firm level.

C0329: A generalization of Benford's law that considers relative quantities

Presenter: **Mario Maggi**, University of Pavia, Italy

Co-authors: Roy Cerqueti, Alex Ely Kossovsky, Claudio Lupi

Benford's law is a discrete probability distribution often observed in the significant digits of many datasets. A wide range of natural phenomena and human-derived data exhibit a high frequency of low significant digits. This behavior is commonly described by Benford's law. Deviations from Benford's law may indicate that the data are bounded, non-numerical (such as codes), or manipulated. This latter point makes Benford's law a useful tool for detecting anomalies and fraud. This paper proposes a generalization of Benford's law, which is named the general law of relative quantities (GLORQ). The GLORQ is a two-parameter discrete probability distribution that is independent of the numeral system used. Like Benford's law, which is a special case, the GLORQ considers data bins that repeat and expand over different orders of magnitude. After introducing the theoretical derivation of the GLORQ, some relations between it and power law variables are presented. Some empirical examples are also provided. Moreover, some simulation studies may suggest the usefulness of the GLORQ in detecting some data manipulations that comply with Benford's law.

C0509: From red flags to red alerts: Prioritizing irregularities in public procurement through predictive modeling

Presenter: **Ivan Pastor**, Universidad Internacional de La Rioja, Spain

Co-authors: Felix J Lopez Iturriaga

While red flag indicators are commonly used to detect potential integrity risks in public procurement, their widespread presence often limits operational effectiveness. A high volume of flagged contracts can overwhelm oversight capacities and dilute attention from truly anomalous or suspicious cases. A complementary approach is proposed to enhance the practical utility of existing red flag systems. A neural network model is introduced to jointly predict two key outcomes for public contracts: The number of bids received and the post-award integrity score assigned by the open tender platform, a composite risk indicator based on procedural transparency, competition, and potential irregularities. By comparing the models' predictions with actual outcomes, contracts whose results deviate significantly and negatively from expected norms are flagged as high-priority cases for audit or further investigation. The empirical application covers over 200,000 high-value public contracts awarded in Spain between 2022 and 2024. The model achieves strong predictive performance ($R^2 = 71.96\%$ for bids; 59.15% for integrity score), enabling a reliable

benchmark of normal behavior. A small but relevant subset of contracts shows unexpected outcomes, such as receiving a single bid when healthy competition is expected. The approach supports smarter oversight and more efficient allocation of audit resources.

C0390: A two-part income-based model for estimating underreported income in the presence of consumption measurement error

Presenter: **Maria Felice Arezzo**, Sapienza University of Rome, Italy

Co-authors: Giuseppina Guagnano, Domenico Vitale, Roberta Di Stefano, Alberto Arcagni

A novel two-part income-based model is introduced to estimate underreported income, addressing limitations of expenditure-based approaches in datasets with consumption measurement error, such as the Bank of Italy's SHIW. Traditional methods often require controlled data or assumptions about evading units, hindering broader applicability. While prior works use expenditure frameworks, they critically assume error-free consumption data, a condition often violated. The model circumvents this by directly incorporating measurement error in reported income. Part one establishes a preliminary classification of statistical units into potential "evader" and "non-evader" categories based on the income-consumption gap, explicitly acknowledging misclassification risk. Part two focuses on the provisional evader group, statistically estimating their expected true income to quantify underreporting, using advanced econometrics to mitigate misclassification and measurement error bias. A significant methodological advancement is offered for estimating income underreporting when comprehensive, error-free consumption data is unavailable, contributing to more robust empirical tools for economic measurement and informing fiscal policy.

CO049 Room BCB 211 CFE SESSION: A TRIBUTE TO H. PESARAN III

Chair: Peter Moffatt

C0393: Feeding inflation: The non-linear spillovers of global food commodities

Presenter: **Stephane Dees**, Banque de France, France

The role of global food price shocks in driving inflation dynamics is examined, with a particular focus on non-linear effects and heterogeneity across countries. Using a Bayesian vector autoregression (BVAR) approach, global food price shocks are identified, and their impact is assessed on inflation. Quantile regressions reveal that inflation's sensitivity to food price shocks is significantly higher in high-inflation regimes, with elasticity estimates being significantly larger for higher percentiles compared to the median. Additionally, panel estimations indicate that this sensitivity is particularly pronounced in emerging markets and developing economies, where food constitutes a larger share of household consumption. Findings contribute to the literature on global inflation synchronization and highlight the importance of accounting for inflation regimes when analyzing the transmission of commodity price shocks.

C0925: Time-varying global VARs with application to interconnectedness, structural analysis, and nowcasting

Presenter: **L Vanessa Smith**, University of York, United Kingdom

Co-authors: Francesco Meglioli

A time-varying parameter global vector autoregressive (TVP-GVAR) model is developed that leverages Kalman filtering techniques with forgetting factors to address the computational challenges typically associated with large-scale time-varying systems. The flexibility and broad applicability of the modelling approach are demonstrated through three distinct empirical applications, highlighting its practical value for regulators and policymakers. First, a baseline version of the TVP-GVAR model is employed, and is used to examine the evolution of interconnectedness between European banks and non-bank financial institutions over the past two decades. Second, the baseline model is extended to enable structural analysis, applying it to investigate the dynamic evolution of government spending multipliers across countries. Finally, the forgetting factors are incorporated within an unscented Kalman filter framework to estimate a MIDAS TVP-GVAR model, making the approach computationally feasible for nowcasting purposes. High-frequency macroeconomic indicators, including industrial production and prices, are used to produce timely nowcasts of quarterly economic activity.

C1277: Distributional effects with two-sided measurement error: An application to intergenerational income mobility

Presenter: **Tong Li**, Vanderbilt University, United States

Co-authors: Brantly Callaway, Emmanuel Tsyawo, Irina Murtazashvili

The aim is to consider the identification and estimation of distributional effect parameters that depend on the joint distribution of an outcome and another variable of interest ("treatment") in a setting with "two-sided" measurement error, that is, where both variables are possibly measured with error. Examples of these parameters in the context of intergenerational income mobility include transition matrices, rank-rank correlations, and the poverty rate of children as a function of their parents' income, among others. Building on recent work on quantile regression (QR) with measurement error in the outcome (particularly, a prior study), it is shown that, given (i) two linear QR models separately for the outcome and treatment conditional on other observed covariates and (ii) assumptions about the measurement error for each variable, one can recover the joint distribution of the outcome and the treatment. Besides these conditions, the approach does not require an instrument, repeated measurements, or distributional assumptions about the measurement error. Using recent data from the 1997 National Longitudinal Study of Youth, it is found that accounting for measurement error notably reduces several estimates of intergenerational mobility parameters.

C0928: A direct test of the fundamental assumption of option pricing models

Presenter: **Peter Moffatt**, University of East Anglia, United Kingdom

Co-authors: Kensley Blaise, Anthony Rezitis

The most fundamental assumption underlying option pricing models is that the market price of an option is equal to the discounted expected value of the final payoff. This assumption is tested directly with data on market prices of options combined with data on realised final payoffs. The data set contains around 1.5 million European options written on the S&P500 index, between January 2022 and December 2023, with expiry dates ranging from January 2022 to July 2025. Only near-the-money options are included in the sample. The framework for testing the hypothesis of interest is a heteroscedastic regression model with discounted actual payoff at expiry as the dependent variable, and market price of the option as the independent variable. The joint hypothesis under test is essentially that the intercept is zero and the slope is one. This hypothesis is tested both parametrically and non-parametrically. As well as being tested for the entire sample, it is tested separately for call options and put options. It is also tested separately for bull-market phases and bear-market phases. The fundamental assumption is always strongly rejected, usually with the intercept being different from zero but the slope not being different from one. It is concluded that a useful way of judging the performance of an option pricing model is to compare computed option valuations to discounted final payoffs, rather than to market prices.

CO134 Room BCB 212 ADVANCES IN TIME SERIES MODELING

Chair: Jui-Chung Yang

C0490: On the consistency of Bayesian adaptive testing under the Rasch model

Presenter: **Yu-Chang Chen**, National Taiwan University, Taiwan

Co-authors: Hau-Hung Yang, Chia-Min Wei

The consistency of Bayesian adaptive testing methods is established under the Rasch model, addressing a gap in the literature on their large-sample guarantees. Although Bayesian approaches are recognized for their finite-sample performance and capability to circumvent issues such as the cold-start problem. However, rigorous proofs of their asymptotic properties, particularly in non-i.i.d. structures, remain lacking. Conditions are derived under which the posterior distributions of latent traits converge to the true values for a sequence of given items, and demonstrate that Bayesian estimators remain robust under the mis-specification of the prior. The analysis then extends to adaptive item selection methods in which items are chosen endogenously during the test. Additionally, a Bayesian decision-theoretical framework is developed for the item selection problem, and

a novel selection is proposed that aligns the test process with optimal estimator performance. These theoretical results provide a foundation for Bayesian methods in adaptive testing, complementing prior evidence of their finite-sample advantages.

C0481: Bayesian analysis of long memory and roughness in financial volatility

Presenter: Toshiaki Watanabe, Hitotsubashi University, Japan

Co-authors: Jouchi Nakajima

Realized volatility (RV) calculated using intraday returns has recently been used as an accurate estimator of financial volatility. Some researchers have documented that the log-RV may follow a long-memory process, which is represented by a fractional Brownian motion with the Hurst exponent greater than 0.5 or a fractionally integrated process with a positive difference parameter. Recent studies show that the log difference in RV may be rough, which is represented by a fractional Brownian motion with the Hurst exponent less than 0.5 or a fractionally integrated process with a negative difference parameter. A discrete-time model that is consistent with these two phenomena is presented, and a Bayesian method is developed for the analysis of this model using Markov chain Monte Carlo. Empirical analysis using the RV of the S&P500 and Nikkei 225 stock indices reveals that the model that takes into account both long-term memory and roughness has the highest volatility prediction accuracy, surpassing the heterogeneous autoregressive (HAR) model, which is known to have high volatility prediction accuracy.

C0501: Approximate maximum likelihood estimation of a threshold diffusion process with separate drift and diffusion regimes

Presenter: Henghsiu Tsai, Academia Sinica, Taiwan

Co-authors: Edward MH Lin

A novel threshold diffusion model in which the drift and diffusion components undergo regime shifts at distinct threshold values is proposed. Specifically, the drift term changes regimes at one threshold, while the diffusion term transitions at a potentially different threshold level. This doubly thresholded specification offers greater flexibility in capturing asymmetric dynamics in both the conditional mean and volatility. To estimate the model, an approximate maximum likelihood estimation (AMLE) procedure is developed that remains computationally tractable despite the model's structural complexity. Monte Carlo simulations confirm the consistency and efficiency of the proposed AMLE method and demonstrate that standard information criteria (AIC, BIC, HQIC) can reliably distinguish between models with shared versus distinct thresholds. Empirical applications using U.S. Treasury interest rates and the 10-Year/3-Month yield spread further support the proposed framework, revealing structural asymmetries that are better captured by allowing separate thresholds in the drift and diffusion components.

C0526: Synthetic historical control for policy evaluation

Presenter: Jui-Chung Yang, National Taiwan University, Taiwan

Co-authors: Yi-Ting Chen, Tzu-Ting Yang

A synthetic historical control method is proposed for policy evaluation without relying on cross-sectional untreated units. The approach builds upon a semi-parametric time-series regression and adapts the conventional synthetic control method by replacing cross-sectional untreated units with historical units. It is demonstrated that the proposed method is approximately unbiased for estimating the intervention effects. Additionally, its effectiveness is shown through both simulation studies and real-world empirical applications. This method offers a valuable alternative for estimating causal effects in scenarios where suitable cross-sectional control units are unavailable.

CO084 Room BCB 213 THEORY AND COMPUTATION FOR STOCHASTIC PROCESS MODELS

Chair: Hiroki Masuda

C0299: Statistical inference for highly correlated stationary point processes and noisy bivariate Neyman-Scott processes

Presenter: Takaaki Shiotani, Graduate School of Mathematical Sciences, The University of Tokyo, Japan

Motivated by the task of estimating lead-lag relationships in high-frequency financial data, the focus is on how to infer highly correlated structures between point processes using quasi-likelihood methods. High correlation refers to situations in which the correlogram diverges for certain parameter values, placing the problem outside the reach of existing asymptotic theory. An asymptotic framework is developed that is applicable under such conditions. As a concrete example, a noisy bivariate Neyman-Scott point process that can be employed for lead-lag estimation is treated.

C0312: Modeling and estimating asset price jumps using CARMA-Hawkes processes and high-frequency VIX data

Presenter: Lorenzo Mercuri, University of Milan, Italy

A self-exciting point process with a continuous-time autoregressive moving average intensity is introduced, referred to as the CARMA(p,q)-Hawkes model. This framework generalizes the classical Hawkes process by replacing the Ornstein Uhlenbeck intensity with a CARMA(p,q) process, where the associated state is driven by the counting process itself. While maintaining the analytical tractability of the Hawkes process, the proposed model captures more intricate temporal structures commonly observed in financial market data. Building on this foundation, a novel asset price model is developed based on a compound CARMA(p,q)-Hawkes process with random jump sizes. This construction can serve as a building block for pure-jump, stochastic volatility jump-diffusion (SVJ) and stochastic volatility model with jumps in the stock price and the volatility (SVJJ). A calibration approach is also proposed, formulated as a maximum likelihood estimation problem. A key challenge in this procedure is the filtering of the unobservable volatility and CARMA-driven intensity processes. While previous studies have relied on the extended Kalman filter or Bayesian filtering techniques, an alternative strategy is explored that recovers these latent processes directly from high-frequency VIX data. Finally, the performance of the estimation method is assessed through both simulation studies and empirical analysis using real data.

C0474: Parameter estimation for weakly interacting hypoelliptic diffusions

Presenter: Yuga Iguchi, Lancaster University, United Kingdom

Co-authors: Alexandros Beskos, Greg Pavliotis

Parameter estimation is studied for an N -weakly interacting particle system (IPS) of multidimensional hypo-elliptic SDEs, where each particle is typically characterized by a degenerate diffusion matrix. A locally Gaussian approximation is proposed, carefully designed to address the degenerate structure, which provides an approximate joint transition density and thus forms a tractable full likelihood, enabling statistical inference for a wider class of hypo-elliptic IPSs that are not covered by a recent work relying on a locally degenerate approximation, specifically the Euler-Maruyama method. A contrast estimator is then analyzed, based on the likelihood from n discretely sampled particle observations with a fixed time interval T , and its asymptotic normality is shown as $n, N \rightarrow \infty$ with a requirement on the step-size $\Delta_n \equiv T/n = o(1/N)$, assuming all particle coordinates are observed. In practical situations where only partial coordinates of particle trajectories are observed, the proposed locally Gaussian approximation offers greater flexibility in inference combined with standard computational statistics methodologies, compared to the Euler-Maruyama type estimator that requires a particular structure for the hypo-elliptic model. Numerical experiments that illustrate the effectiveness of using the locally Gaussian approximation-based likelihood are presented in settings where complete or partial particle trajectories are observed.

C0494: Lead-lag analysis of high-frequency financial data based on point processes

Presenter: Yuta Koike, University of Tokyo, Japan

Co-authors: Takaki Hayashi, Takaaki Shiotani

A new theoretical framework is developed to analyze lead-lag relationships between the order arrivals of two assets. A seminal work proposed model-free measurements of cross-market trading activity based on cross-counts of order arrivals and demonstrated that they can sharply identify their lead-lag relationships, but their mathematical meanings remained unclear. To resolve this issue, the problem of estimating lead-lag relationships in order is formulated as a problem of estimating the shape of the so-called cross-pair correlation function (CPCF) of a bivariate stationary

point process, which has been extensively studied in the literature of biostatistics. Then, the lead-lag time parameter is defined as its maximizer. Within this framework, the peak in the prior study's cross-market activity measure can be interpreted as an estimator for the lead-lag time parameter. Furthermore, an alternative lead-lag time estimator is proposed based on kernel density estimation, and it is shown that it has desirable theoretical properties along with superior numerical performance.

CO238 Room BCB M202 ROBUST METHODS IN ENERGY, ENVIRONMENT, AND FINANCE
Chair: Luigi Grossi
C0316: Electricity price forecasting with LSTM: Evidence from the Italian day-ahead power market

Presenter: **Filippo Beltrami**, Eurac Research, Italy

Co-authors: Luigi Grossi, Patrick Thoeni, Riccardo Lucarno

The aim is to present a comprehensive forecasting framework for electricity prices in the Italian day-ahead power market, leveraging both traditional time-series models and advanced artificial intelligence (AI) techniques. The focus is on a novel, high-resolution proprietary dataset which has not yet been used in the existing literature, spanning from 2021 to 2024 at 15-minute intervals and encompassing, among other variables, market volumes, zonal prices, detailed weather metrics (e.g., temperature), cross-border exchanges with Switzerland, France, and Austria, as well as plant maintenance schedules. The modeling pipeline systematically implements a suite of long short-term memory (LSTM) architectures on this uniquely rich and granular dataset. Each model is rigorously benchmarked against established econometric and machine-learning references, namely a refined lasso-estimated autoregressive (LEAR) model and a two-layer DNN. Model performance is evaluated using a comprehensive set of metrics (rMAE, MAE, MAPE, sMAPE, RMSE) alongside measures of computational effort, with the goal of identifying the optimal trade-off between forecasting accuracy and computational cost and thus aligning cutting-edge methodological advances with the practical performance requirements of industry stakeholders.

C0447: A robust empirical analysis of the EU carbon border adjustment mechanism: Early signals from EU and India steel trade

Presenter: **Gian Luca Vriz**, University of Edinburgh, United Kingdom

Co-authors: Theodor Cojoianu, Carolyn Fischer, Luca Taschini

The European Union's carbon border adjustment mechanism (CBAM) aims to reduce carbon leakage by applying a carbon price on imports from countries with less stringent environmental standards. The initial impact of CBAM is examined on the iron and steel trade between the European Union and India. Using firm-level emissions data and trade data, a robust empirical framework is developed to classify a sample of Indian exporters by emission intensity and track changes in their export behavior during the CBAM's reporting phase. A difference-in-differences methodology is applied to assess variations in shipment sizes and unit prices. The results suggest that CBAM influences firm competitiveness and supply chain configurations, even before its full implementation in 2026. The analysis emphasizes the importance of policymakers to consider carbon intensity in their international trade strategies and highlights the value of firm-level data in assessing the real-world impacts of environmental policies.

C0839: The hard, transition, and the good times

Presenter: **Emmanuel Fianu**, De Montfort University, United Kingdom

Co-authors: Daniel Felix Ahelegbey, Roberto Casarin, Luigi Grossi

The aim is to expand the analysis of various interactions spanning over two decades to examine the extreme downside interconnectedness between crude oil and agricultural commodity futures markets. In effect, it captures and provides updated information that details the degree of interrelatedness of various unique abnormal market situations. By focusing on tail dependencies, the aim is to uncover the varying structural dynamics that drive co-movements during periods of market stress. Utilizing daily price data for crude oil, ten agricultural commodities, and the Bloomberg Commodity Index, advanced econometric models are employed to quantify the extreme downside risk and their interactions via a novel, innovative network structure. It is also investigated which markets act as a flight to safety during extreme events, and it is identified that markets that dominate risk transmission or are vulnerable to tail risk propagation. Results show that soya bean meal and soya bean are the major recipients of tail risk in all the periods on average, whilst corn is a major transmitter of tail shock that permeates through all the various scenarios. The empirical evidence suggests varying degrees of interconnection among the commodity market index, oil commodity futures, and agricultural commodity futures markets, coupled with a decreasing trend of hedging effectiveness.

C0860: Robust estimation of mixed models for energy production

Presenter: **Fabrizio Laurini**, University of Parma, Italy

Co-authors: Luis Angel Garcia-Escudero, Agustin Mayo-Isacar

Parameter estimation for mixed models can be severely biased by outliers. The outliers can be in the conditional distribution of the response variable or even in the explanatory variables. The proposal is to mitigate the effect of outliers by introducing robust estimations by trimming a fixed proportion of observations. The level of trimming is an input hyperparameter, and it can be adjusted and monitored with proper diagnostics. Some results are shown with artificial data, and it is compared with available competitors. The proposed robust method is also applied to real data, based on electricity production data in Europe.

CO027 Room BCB 308 STATISTICAL LEARNING FOR COMPLEX AND DYNAMIC SYSTEMS
Chair: Matteo Borrotti
C0926: SAFE: A unified framework for the validation of tabular, imaging and longitudinal synthetic data in clinical research

Presenter: **Gianluca Asti**, IRCCS Humanitas, Italy

Co-authors: Elena Zazzetti, Elisabetta Sauta, Mattia Delleani, Eleonora Iascone, Alessandro Bruseghini, Giulia Maggioni, Luca Lanino, Alessia Campagna, Marta Ubezio, Alessandro Buizza, Gabriele Todisco, Cristina Astrid Tentori, Antonio Russo, Alessandra Crespi, Nicole Pinocchio, Maria Chiara Grondelli, Alessandro Forcina Barrero, Viktor Savevski, Armando Santoro, Saverio D'Amico, Matteo Giovanni Della Porta

SAFE (Synthetic vAlidation FramEwork) is a statistically grounded, scalable framework designed to rigorously validate synthetic data (SD) across clinical modalities, including structured, temporal, and imaging data. It focuses on three core characteristics: Fidelity, utility, and privacy, through a suite of quantitative metrics tailored to data type and clinical context. For tabular and longitudinal data, SAFE applies RMSE, R, total variation distance (TVD), Kolmogorov-Smirnov tests, SMAPE, and correlation analysis to assess distributional alignment. For imaging, fidelity is evaluated using FID and MS-SSIM. Privacy is quantified via membership inference attacks, while clinical utility is tested through downstream tasks such as survival prediction (L1-penalized Cox models), disease classification (e.g., XGBoost), and synthetic control arm construction. SAFE was validated in two domains: Synthetic bone marrow images for hematologic malignancies and longitudinal breast cancer records. In the hematology use case, synthetic images matched real histopathological features and improved disease classification by 10% (F1 score) and survival modeling by over 10% (C-index). In breast cancer, LLMs like Mistral-7B generated high-fidelity, privacy-preserving data that enhanced predictive modeling and successfully emulated clinical trial outcomes. SAFE enables statistically robust integration of SD into clinical research, supporting precision medicine and real-world evidence generation.

C0798: Dependent spike-and-slab for drug combinations in diabetic kidney disease

Presenter: **Jiefeng Bi**, University of Milano-Bicocca, Italy

Diabetic kidney disease (DKD) is a serious long-term complication of type II diabetes mellitus, primarily resulting from glucose metabolism disturbances associated with insulin resistance. Bayesian augmented learning (BAL) is employed, which leverages information across stages via dependent spike-and-slab priors to enhance treatment adaptation. False discovery rate (FDR) is also introduced as a method to identify the most significant variables at each intervention stage. This approach is applied to data from the PROVALID study, a prospective observational cohort of

type II diabetes mellitus patients designed for biomarker validation. To adapt BAL for the PROVALID dataset, the original two-stage framework is extended to a four-stage model. Furthermore, a Bayesian predictive model is developed to determine optimal drug combinations for improving prognosis in individual DKD patients. The approach successfully identifies key molecular biomarkers and clinical parameters at each treatment stage, specifically in relation to drug effects on estimated glomerular filtration rate (eGFR).

C0835: Advancing statistical jump models: A novel method for robust temporal clustering with state-dependent feature selection

Presenter: **Federico Cortese**, University of Milan, Italy

Co-authors: Antonio Pievatolo

Statistical sparse jump models offer a flexible alternative to traditional hidden Markov models by capturing complex dynamics via smooth transitions between regimes while simultaneously performing feature selection. However, they typically select features globally, assuming their relevance across all states, and might be sensitive to outliers. The aim is to introduce a robust sparse jump model that estimates a latent state sequence over time, along with a matrix of state-specific feature weights. This allows for modeling regime-dependent feature relevance, which might be crucial in applications where different subsets of features possibly drive different regimes. Robustness is ensured by: (1) integration of an outlyingness factor that down-weights atypical observations during estimation, and (2) the use of medoids, rather than prototypes, as state representatives.

C0826: Neural network methods for volatility forecasting

Presenter: **J Miguel Marin**, University Carlos III, Spain

Co-authors: Helena Veiga, Hongfei Guo

Accurate volatility forecasting is essential for risk management and financial market investment decisions. The aim is to present a new framework that improves volatility prediction over various time horizons by combining neural networks and stochastic volatility models. The method uses information from both low- and high-frequency financial data and provides symmetric and asymmetric stochastic volatility specifications enhanced with jump components. Uncertainty is quantified robustly by estimating model parameters using Bayesian inference techniques. To improve predictive performance and stability, a Bayesian ensemble method is also created that combines forecasts from multiple models. Forecast accuracy is evaluated using two loss functions and conditional superior predictive ability tests when applying the suggested methodology to three significant international stock indices. Findings demonstrate significant improvements in volatility forecasting accuracy, highlighting the advantages of integrating machine learning and traditional methods within a unified Bayesian framework.

CO138 Room BCB 309 MODERN STATISTICAL LEARNING AND INFERENCE FOR STRUCTURED DATA

Chair: Li Ma

C0273: Weight matrix compression based on PDB model in deep neural networks

Presenter: **Zeng Li**, Southern University of Science and Technology, China

Weight matrix compression has been demonstrated to effectively reduce overfitting and improve the generalization performance of deep neural networks. Compression is primarily achieved by filtering out noisy eigenvalues of the weight matrix. A novel population double bulk (PDB) model is proposed to characterize the eigenvalue behavior of the weight matrix, which is more general than the existing population unit bulk (PUB) model. Based on the PDB model and random matrix theory (RMT), a new PDBLS algorithm is discovered for determining the boundary between noisy eigenvalues and information. A PDB noise-filtering algorithm is further introduced to reduce the rank of the weight matrix for compression. Experiments show that the PDB model fits the empirical distribution of eigenvalues of the weight matrix better than the PUB model, and the compressed weight matrices have a lower rank at the same level of test accuracy. In some cases, the compression method can even improve generalization performance when labels contain noise.

C0612: Few-shot personalization for nonparametric regression with minimax optimality

Presenter: **Sai Li**, Renmin University of China, China

Co-authors: Linjun Zhang

Personalized modeling in nonparametric regression aims to adapt pre-trained models to individual-specific data with minimal samples, addressing the crucial challenge of few-shot learning. A theoretical framework is established for few-shot personalization in nonparametric regression. Novel algorithms are proposed that adapt classical nonparametric estimation techniques to the personalized setting. Results provide rigorous guarantees and offer new directions for designing data-efficient, personalized models in real-world applications where data scarcity is a key concern.

C0626: Causal inference in biomedical imaging via functional linear structural equation models

Presenter: **Ting Li**, Shanghai University of Finance and Economics, China

Understanding the causal effects of organ-specific features from medical imaging on clinical outcomes is essential for biomedical research and patient care. A novel functional linear structural equation model (FLSEM) is proposed to capture the relationships among clinical outcomes, functional imaging exposures, and scalar covariates like genetics, sex, and age. Traditional methods struggle with the infinite-dimensional nature of exposures and complex covariates. FLSEM overcomes these challenges by establishing identifiable conditions using scalar instrumental variables. The functional group support detection and root finding (FGS-DAR) algorithm is developed for efficient variable selection, supported by rigorous theoretical guarantees, including selection consistency and accurate parameter estimation. A test statistic is further proposed to test the nullity of the functional coefficient, establishing its null limit distribution. The approach is validated through extensive simulations and applied to UK Biobank data, demonstrating robust performance in detecting causal relationships from medical imaging.

CO037 Room BCB 310 ADVANCES IN FUNCTIONAL AND COMPLEX DATA ANALYSIS

Chair: Alessia Pini

C0789: Should we correct the bias in confidence bands for repeated functional data?

Presenter: **Emilie Devijver**, CNRS, France

Co-authors: Adeline Leclercq-Samson

While confidence intervals for finite quantities are well-established, constructing confidence bands for objects of infinite dimension, such as functions, poses challenges. The concept of parametric confidence bands is explored for functional data with an orthonormal basis. Specifically, the method proposed by a prior study is revisited, which yields confidence bands for the projection of the regression function in a fixed-dimensional space. This approach can introduce bias in the confidence bands when the dimension of the basis is misspecified. Leveraging this insight, a corrected, unbiased confidence band is introduced. Surprisingly, the corrected band tends to be wider than what a naive approach would suggest. To address this, a model selection criterion that allows for data-driven estimation of the basis dimension is proposed. The bias is then automatically corrected after dimension selection. These strategies are illustrated using an extensive simulation study. It is concluded with an application to real data.

C0824: PDE-regularized models for spatiotemporal and quantile regression

Presenter: **Eleonora Arnone**, University of Turin, Italy

Co-authors: Laura Sangalli

The focus is on the development of flexible statistical models for spatiotemporal data, designed to accommodate complex structural features such as anisotropy, non-stationarity, and domain irregularity. It begins by introducing a class of semiparametric regression models that incorporate partial differential equations (PDEs) as a means of regularization. These models employ differential operators in the penalty term to encode structural assumptions or prior scientific insight, such as geometric constraints or directional trends within the estimation procedure. The use

of PDE-based roughness penalties allows for accurate modelling over irregular spatial domains and naturally captures anisotropic behaviors. The framework is then extended to the spatiotemporal quantile regression setting, where the modelling target shifts from the conditional mean to specific distributional quantiles of the response. This is particularly relevant in contexts where the variability and tail behavior of the process are of interest. A spatiotemporal quantile regression model is formulated in which the quantile field is estimated by minimizing a pinball loss function, regularized by space-time differential penalties. The resulting model allows for a nonparametric representation of the quantile function and accommodates complex spatial and temporal dependence.

C1065: A conformal prediction approach to predict populations of networks

Presenter: **Matteo Fontana**, Royal Holloway, University of London, United Kingdom

Co-authors: Anna Calissano, Simone Vantini, Gianluca Zeni

Despite the growing importance of population of network data and its analysis, a significant gap exists in the literature with respect to predicting these objects and quantifying the certainty of those predictions. To address this, a new technique is proposed for generating statistically valid prediction sets for populations of networks. The method, rooted in conformal prediction, uniquely handles both labelled (fixed-node) and unlabeled (variable-node) graph structures, the latter by defining sets in a discrete quotient space. This approach offers three key advantages: It requires no distributional assumptions, provides finite-sample guarantees, and yields results that are easily interpretable. Simulation studies confirm the method's theoretical properties, while an analysis of player passing networks from the FIFA 2018 World Cup showcases its real-world applicability.

C1239: Detecting shape outliers in functional data via complexity-driven mixture modeling

Presenter: **Enea Bongiorno**, Università del Piemonte Orientale, Italy

Co-authors: Kwo Lik Lax Chan, Aldo Goia

The aim is to explore the challenge of identifying atypical shapes within functional data by framing them as the outcome of contamination from high-complexity elements in a mixture model setting. The focus is on three key directions: (i) introducing a complexity measure derived from small ball probabilities, (ii) formulating a mixture model that incorporates this notion of complexity and analyzing its theoretical impact on small ball behavior, and (iii) developing a decomposition algorithm that separates the mixture into interpretable components, thereby facilitating the implicit detection of outliers. The effectiveness of the proposed framework is illustrated through an applied case study.

CO267 Room BCB 311 CLUSTERING OF HETEROGENEOUS DATA

Chair: Efthymios Costa

C0522: A unified approach to outlier identification for mixed type data

Presenter: **Christian Hennig**, University of Bologna, Italy

Co-authors: Efthymios Costa

An approach for identifying outliers in data with continuous as well as ordinal variables is presented with possible extension to nominal categorical variables. The approach is based on robust Mahalanobis distances (based on FastMDC) and a definition of outliers as observations that are in low probability regions relative to a multivariate Gaussian distribution. In order to unify the contribution of continuous and ordinal variables to the robust Mahalanobis distance and the definition of outliers, ordinal variables are modeled as stemming from thresholding a latent Gaussian (one for each ordinal variable). Polychoric and polyserial correlation are used to estimate the covariance matrix of the underlying multivariate Gaussian distribution. There can be issues with singularity, particularly in cases in which a single category of a variable contains a large percentage of observations, which may force FastMCD to estimate its variance as zero. For this reason, the covariance matrix will be regularized. Nominal categorical variables can be incorporated by introducing dummy variables for the categories.

C0905: Variable weighting for mixed-type data clustering

Presenter: **Marianthi Markatou**, University at Buffalo, United States

Co-authors: Alexander Foss

Conceptual and technical challenges are discussed in variable weighting and variable selection in cluster analysis of mixed-type data. The specific objectives of existing weighting procedures are quite different and are often selected without a thorough discussion of competing objectives and consideration of the implications of these selection choices. The different objectives of existing weighting and variable selection techniques are described and investigated, and then the variable weighting technique MEDEA (multivariate eigenvalue decomposition error analysis) is developed. It is shown that MEDEA weighting fills an important and previously unfilled niche in cluster analysis of mixed-type data, as demonstrated by its superior performance in numerous experiments with both synthetic and actual datasets.

C0945: Nonparametric kernel and spline cluster weighted models

Presenter: **John Thompson**, The University of British Columbia, Canada

Co-authors: Ling Xue

Cluster weighted models (CWMs) are a class of finite mixture of regression models that jointly model random covariates and response variables. A challenge in CWMs is specifying the correct parametric functional form of each cluster's data-generating process (DGP), as well as the distribution of covariates. Nonparametric regression estimators, such as those using kernel and spline functions, can be employed in CWMs to estimate unspecified nonlinear regression functions for each cluster. These estimators do not require specification of the functional form of the DGP, but rather satisfy some smoothness and moment conditions. CWMs with nonparametric regression function estimators can be estimated using an expectation-maximization algorithm. Model smoothness during estimation is controlled using a roughness penalty on spline coefficients and cross-validated bandwidth selection for kernel functions. Kernel and spline CWMs are applied to simulated and real datasets with various linear and nonlinear cluster shapes, and their performance is compared to that of state-of-the-art CWM methods.

C1035: Cluster analysis of directional data based on data depths

Presenter: **Giuseppe Pandolfo**, University of Naples Federico II, Italy

Co-authors: Antonio D Ambrosio

A new depth-based clustering procedure for directional data is proposed. Such a method is fully non-parametric and has the advantage of being flexible and applicable even in high dimensions when a suitable notion of depth is adopted. The introduced technique is evaluated through an extensive simulation study. In addition, a real data example in text mining is given to explain its effectiveness in comparison with other existing directional clustering algorithms.

CO126 Room BCB 312 DESIGN OF EXPERIMENTS METHODOLOGY AND APPLICATIONS

Chair: Kalliopi Mylona

C0362: New advances in algorithmic techniques for computing optimal designs

Presenter: **Irene Garcia Camacha Gutierrez**, University of Castilla-La Mancha, Spain

Co-authors: Ricardo Negrete Gallego, Sergio Pozuelo Campos

One of the biggest challenges in calculating optimal experimental designs is coping with the computational cost. Algorithmic techniques remain key in this field: While analytical solutions are often impractical, numerical techniques have become the most effective option. An innovative algorithm that combines the algorithms traditionally used in optimal experimental design with metaheuristic techniques is provided. For this purpose, the foundations of Wynn-Fedorov, multiplicative, and the particle Swarm optimization algorithms are suitably adapted. This combination not only enhances computational efficiency but also is more robust in facing complex optimization problems.

C0419: Constructing $E(s^2)$ -optimal two-level supersaturated designs via orthogonal arrays*Presenter:* **Emmanouil Androulakis**, University of Piraeus, Greece*Co-authors:* Haralambos Evangelaras, Kashinath Chatterjee

The purpose is to investigate $E(s^2)$ -optimality of two-level supersaturated designs, which are constructed by augmenting an orthogonal array with two-column interaction terms. By revisiting Wu's construction method, it is shown that $E(s^2)$ -optimal supersaturated designs can be obtained when the starting design is a two-level orthogonal array with n runs and $n-1$, $n-2$, or $n-3$ columns, as long as its two-column interactions are partially aliased. Results rely on the framework of the generalized wordlength pattern and its connection to the concept of J-characteristics, offering new insights into the structure and performance of supersaturated designs derived from orthogonal arrays.

C0430: Split-plot experiments with replicated runs in pharmaceutical synthesis*Presenter:* **Martin Otava**, Johnson & Johnson, Czech Republic*Co-authors:* Kalliopi Mylona

Restricted randomized designs are essential in the pharmaceutical synthesis due to the operational restrictions and their cost-effectiveness compared to entirely randomized designs. Specifically, the split-plot designs are very effective in reducing the cost of an experiment in the presence of hard-to-change factors and/or of multi-stage processes. In classical designs, replicated runs for pure-error estimation are commonly employed, but they are rarely used in the more complex setting of restricted randomization. The reason is that, in practice, experiments in industry can rarely follow all the assumptions/conditions that are included in the methodological papers. A split-plot design based on a Bayesian D-optimality criterion can be adapted to ensure more precise pure-error estimation of the variance components and to fit the additional needs, the speedy implementation, and the restrictions that a real case scenario in industry often imposes. Emphasis is placed on the practical aspect of how to modify the complex design in a way that keeps the desired qualities and on how to assess the impact of more arbitrary changes.

C0465: Response surface designs for general crossed and nested multi-stratum structures*Presenter:* **Steven Gilmour**, KCL, United Kingdom*Co-authors:* Luzia Trinca

Response surface designs are usually described as being run under complete randomization of the treatment combinations to the experimental units. In practice, however, it is often necessary or beneficial to run them under some kind of restriction to the randomization, leading to multi-stratum designs. In particular, some factors are often hard to set, so they cannot have their levels reset for each experimental unit. A general solution is presented for designing response surface experiments in any multi-stratum structure made up of crossing and/or nesting of unit factors. A stratum-by-stratum approach to constructing designs using compound optimal design criteria is used and illustrated. It is shown that good designs can be found even for large experiments in complex structures.

CO174 Room BCB 313 SURVIVAL ANALYSIS METHODS FOR PARTLY INTERVAL-CENSORED DATA**Chair: Jun Ma****C0983: A penalized likelihood approach for joint modelling of longitudinal covariates and partly interval-censored data***Presenter:* **Annabel Webb**, Cerebral Palsy Alliance Research Institute, Australia*Co-authors:* Jun Ma

The focus is on the joint modeling of longitudinal covariates and partly-interval censored time-to-event data. Longitudinal time-varying covariates play a crucial role in obtaining accurate clinically relevant predictions using a survival regression model. However, these covariates are often measured at limited time points and may be subject to measurement error. Further methodological challenges arise from the fact that, in many clinical studies, the event times of interest are interval-censored. A model that simultaneously accounts for all these factors is expected to improve the accuracy of survival model estimations and predictions. Joint models that combine longitudinal time-varying covariates with the Cox model are considered for time-to-event data, which is subject to interval censoring. The proposed model employs a novel penalized likelihood approach for estimating all parameters, including the random effects. The covariance matrix of the estimated parameters can be obtained from the penalized log-likelihood. The performance of the model is compared to an existing method under various scenarios. The simulation results demonstrated that the new method can provide reliable inferences when dealing with interval-censored data. Data from the Anti-PD1 brain collaboration clinical trial in advanced melanoma is used to illustrate the application of the new method.

C0988: Competing risks analysis using mixture cure cause-specific hazard models with partly interval censoring*Presenter:* **Serigne Lo**, University of Sydney / Melanoma Institute Australia, Australia*Co-authors:* Joseph Descallar, Houying Zhu, Jun Ma

In medical studies, competing risks survival data are commonly estimated using partial likelihood methods, with the assumption that any right-censored patient will experience one of the competing events beyond the conclusion of the study period. In some cases, however, a patient may be considered "cured" from the risks of interest, meaning that none of the risks will occur, resulting in a cured fraction. Furthermore, if disease progression is the event of interest, the exact event times are unknown and often subject to interval censoring. When dealing with such data, employing the standard analysis approach based on cause-specific hazards models while ignoring the cured fraction could lead to biased parameter estimates. A novel approach is introduced that accommodates interval censoring and a cured fraction within cause-specific Cox models when analyzing competing risks data. More specifically, a new maximum penalized likelihood approach is proposed to simultaneously estimate logistic regression parameters for the cured fraction, cause-specific Cox models regression coefficients, and their non-parametric baseline hazards. Asymptotic properties are developed, and simulation studies show reduced bias and improved coverage probability compared with the partial likelihood approach with a mid-point imputation.

C1000: Stable likelihood-based parameter estimation for semi-parametric accelerated failure time models*Presenter:* **Aishwarya Bhaskaran**, Macquarie University, Australia*Co-authors:* Jun Ma

Cox proportional hazards and accelerated failure time (AFT) models are two principal frameworks for assessing covariate effects in survival analysis. The AFT model is particularly appealing as it offers a straightforward interpretation of how covariates impact the event time. While semiparametric AFT models allow for flexible modelling without specifying the baseline hazard, they pose significant computational challenges. Unlike Cox models, estimation in semiparametric AFT models is markedly more demanding, as the transformation of failure times used to estimate the baseline hazard is dependent on the regression coefficients and must be updated iteratively, often leading to instability or convergence issues. Moreover, incorporating interval censoring complicates the likelihood expression, introducing algebraic challenges in deriving gradients and Hessians, particularly when the intervals are narrow. To address these challenges, a maximum penalized likelihood approach is proposed using rescaling techniques for fitting semiparametric AFT models to data with partly interval-censored failure times. The method employs M-splines to approximate the nonparametric baseline hazard, a constrained optimization framework for stable coefficient estimation, and an automatic smoothing selection criterion that accounts for active constraints. Simulation studies demonstrate the robustness and accuracy of our approach across a range of censoring scenarios.

C1074: A new EM algorithm for Cox models with partly interval censoring*Presenter:* **Jun Ma**, Macquarie University, Australia

Partly interval censoring is a common feature in medical datasets. Existing estimation methods for semi-parametric Cox models with partly

interval-censored event times include EM algorithms and direct constrained optimization techniques. A new EM approach tailored to this context is presented. The method treats left- and interval-censored data as missing event times. In each iteration, the E-step computes the conditional expectation of the complete-data log-likelihood. In the M-step, a profile likelihood is used to efficiently update the regression coefficients, followed by an update of the baseline hazard function, which is approximated using a piecewise constant function. This iterative procedure converges fast. A simulation study is presented, along with a real-data example to illustrate the method.

CO275 Room BCB 402 SPATIO-TEMPORAL MODELLING FOR ENVIRONMENTAL AND CLIMATE APPLICATIONS Chair: Ana C Cebrian

C0965: INLA-based Bayesian spatiotemporal models for downscaling hurricane wind speeds on the U.S. north Atlantic coast

Presenter: **Concepcion Ausin**, Universidad Carlos III de Madrid, Spain

Co-authors: Michael Wiper, Ali Sarhadi

The aim is to develop a Bayesian spatiotemporal modeling framework to assess wind-related risk and return periods of major tropical cyclones along the U.S. North Atlantic coast, accounting for nonstationary and warming climate conditions. The analysis is based on a high-resolution downscaled dataset derived from a large ensemble of synthetic storm tracks, generated to be consistent with large-scale circulation patterns from climate models. A binomial generalized linear model is employed, and the probability of extreme wind events is estimated at each location over time using the integrated nested Laplace approximation (INLA), which offers substantial computational efficiency over traditional MCMC methods. Both parametric and semi-parametric structures are incorporated to capture spatial and temporal variations in wind intensity distributions, enabling robust projections of future wind-related hazards across coastal sites.

C0478: Spatiotemporal data fusion for environmental applications

Presenter: **Craig Wilkie**, University of Glasgow, United Kingdom

Co-authors: Claire Miller, Marian Scott, Surajit Ray, Daniela Castro-Camilo, Daniela Cuba, Stephen Jun Villejo, Pietro Colombo

A summary of recent work is presented on data fusion methods for spatiotemporal environmental applications. The increasing availability of environmental data from multiple sources, such as satellites and low-cost sensors, provides an improved understanding of the changing environment. Data from these sources can, however, be of varying quality, often on different spatial and temporal scales. Data fusion approaches aim to combine information from multiple complementary data sources to provide an enhanced understanding of environmental variables compared to using any single data source, with associated uncertainty measures accounting for differences in the quality of information from each source. A summary of recent work is presented on data fusion methods, with a focus on some work by early-career researchers.

C0631: Developing a spatiotemporal NO2 model to assess extreme exposures in the Rome region

Presenter: **Edoardo Rosci**, Sapienza University of Rome, Italy

Co-authors: Jorge Castillo-Mateo, Jesus Asin, Ana C Cebrian, Giovanna Jona Lasinio, Massimo Stafoggia

Extreme air pollution is a major risk factor for population health, and reducing its frequency and intensity has proven benefits for quality of life and life expectancy. To this end, developing updated and robust quantitative methods that can capture the impact of human behavior is essential. Therefore, this can provide public decision-makers with tools to allocate health resources based on health risks. Extreme NO2 concentrations are analyzed in the Lazio region between 2017 and 2022. Following a data-driven approach, key temporal patterns are first explored, such as seasonality, workday effects, and autoregressive structure. Local quantile regression models are then fitted at each monitoring station to capture station-specific behaviors and spatial heterogeneity. Model selection was implemented taking advantage of the quantile models' specific metrics. Finally, the latest developments toward a comprehensive Bayesian hierarchical spatiotemporal model are presented. Future work includes scaling the model to the national level, integrating preferential sampling within the MCMC estimation, and leveraging the inverse chemical relationship between NO2 and O3 to develop multivariate quantile regression models.

C0285: Spatiotemporal modeling of record-breaking daily maximum and minimum temperatures

Presenter: **Jorge Castillo-Mateo**, University of Zaragoza, Spain

Co-authors: Zeus Gracia, Jesus Asin, Ana C Cebrian, Alan Gelfand

The increasing occurrence of record-breaking temperatures presents a critical global challenge. To investigate the influence of climate change on these extreme events, daily maximum and minimum temperature data from 40 meteorological stations across peninsular Spain were analyzed over the period 1960-2023. The dataset is encoded as a daily pair of binary events, indicators, for that day, of whether a yearly upper record was broken for the daily maximum temperature and/or for the daily minimum temperature. A Bayesian hierarchical spatiotemporal bivariate probit regression model was developed to study the probability of such events. The model incorporates an explicit long-term trend, vector autoregressive components, spatial effects based on distance to the coast, relevant interactions, and strong daily spatial random effects. The spatial random effects are modeled using coregionalized Gaussian processes, where covariates are incorporated both in the coregionalization structure and in the covariance function to account for non-stationarity and anisotropy. Model-based inference reveals a strong correlation between record-breaking maxima and minima, with distinct spatial and temporal patterns of climate change signals, as well as differing degrees of persistence and spatial dependence in the two processes.

CO085 Room BCB 403 ADVANCES IN STATISTICS AND ML FOR BREEDING STUDIES

Chair: Pariya Behrouzi

C0305: Regression approaches for modelling genotype-environment interaction and making predictions into unseen environments

Presenter: **Maksym Hrachov**, University of Hohenheim, Germany

Co-authors: Hans-Peter Piepho, Niaz Md Farhat Rahman, Waqas Ahmed Malik

In plant breeding and variety testing, there is an increasing interest in making use of environmental information to enhance predictions for new environments. Linear mixed models that have been proposed for this purpose are reviewed, with an emphasis on predictions and on methods to assess the uncertainty of predictions for new environments. The point of departure is straight-line regression, which may be extended to multiple environmental covariates and genotype-specific responses. When observable environmental covariates are used, this is also known as factorial regression. Early work along these lines can be traced back to Finlay-Wilkinson regression dating back to 1930s. This method, in turn, has close ties with regression on latent environmental covariates and factor-analytic variance-covariance structures for genotype-environment interaction. Extensions of these approaches - reduced rank regression, kernel- or kinship-based approaches, random coefficient regression, and extended Finlay-Wilkinson regression - are the focus of comparison. The objective is to demonstrate how seemingly disparate methods are very closely linked and fall within a common model-based prediction framework. Options are considered for assessing uncertainty of predictions, including cross-validation and model-based estimates of uncertainty, and tested methods on the long-term rice variety trial dataset.

C0377: Evaluating the impact of trait measurement error on genetic analysis of computer vision-based phenotypes

Presenter: **Gota Morota**, The University of Tokyo, Japan

Quantitative genetic analysis of image-derived phenotypes is increasingly being performed for a wide range of traits. Pig body weight estimated by a conventional approach or a computer vision system can be considered as two different measurements of the same trait, but with different sources of phenotyping error. Previous studies have shown that trait measurement error, defined as the difference between manually collected phenotypes and image-derived phenotypes, can be influenced by genetics, suggesting that the error is systematic rather than random and is more likely to lead to misleading quantitative genetic analysis results. The effect of trait measurement error is investigated on the genetic analysis of pig body weight (BW). Genomic heritability estimates for trait measurement error were consistently negligible, regardless of the choice of computer vision

algorithm. In addition, genome-wide association analysis revealed no overlap between the top markers identified for scale-based BW and those associated with trait measurement error. No evidence is found that the BW trait measurement error could be influenced by genetic factors. This suggests that trait measurement error in pig BW does not contain systematic errors that could bias downstream genetic analysis.

C0383: Genomic prediction: A robustness comparison of machine learning approaches

Presenter: Vanda Lourenco, NOVA University of Lisbon and NOVA.id.FCT, Portugal

Co-authors: Joseph O. Ogutu, Piepho Hans-Peter

Accurate estimation of genomic breeding values underpins effective genomic selection in plants and animals. Genomic prediction leverages dense SNP markers and demands models capable of handling extreme dimensionality; machine-learning (ML) algorithms are natural candidates. While many studies benchmark individual ML algorithms, comparisons across algorithmic families remain scarce, and even fewer explore how data contamination affects performance. Yet, breeders routinely confront noisy phenotypes, making robustness as important as raw accuracy. This gap is filled by comparing three supervised ML families: Regularized regression, ensemble learners, and instance-based methods under pristine and contaminated conditions. Using a simulated animal-breeding population, escalating proportions of contaminated phenotypes are imposed, then predictive accuracy (PA) and mean-squared error are quantified. PA declines and MAPE rises with greater contamination and mean shifts; radial outliers impair predictions more than point-mass. Random forest, GLasso, SGB, aENET, and SVM run markedly slower. Results illuminate trade-offs among speed and robustness, and reveal circumstances where penalized regressions outperform more complex alternatives. Guidance is provided for breeders selecting algorithms when data quality is uncertain, emphasizing the need to match model choice to anticipated contamination rather than relying solely on headline accuracy in idealized datasets used for benchmarking.

C0742: Statistical and ML options for prediction of GxE in plant breeding

Presenter: Fred van Eeuwijk, Wageningen University, Netherlands

In plant breeding and genetics, a central research objective is the prediction of phenotypic performance from genetic and environmental inputs. In its simplest form, the question is how to predict single or multiple phenotypes (responses) from genetic and environmental factors and covariates. The default model class for such predictions was that of linear mixed models. Over the last decades, genetic and environmental inputs have become high-dimensional, and additional classes of inputs, like phenomics, are accessible via new types of sensors and measurement devices. Environmental and phenomic inputs are often longitudinal. The mixed model framework for prediction of phenotypes has been extended to incorporate multiple sets of high-dimensional inputs via multiple kernels. The longitudinal aspect of environmental and phenomic data can be addressed by the insertion of various types of base functions or special types of variance-covariance matrices. Alternative methods for analysis and prediction of phenotypes try to exploit the high resolution of some of the environmental and phenomic inputs. Hierarchical Bayesian approaches that model longitudinal phenotypes by systems of differential equations are found. Machine learning and deep learning approaches have been proposed, too. The prediction of phenotypes is considered via different modeling classes with special attention, for the way in which genotype by environment interactions are addressed.

CO320 Room BCB 405 TOPICS IN BAYESIAN MODELLING AND COMPUTATION

Chair: Khue-Dung Dang

C0291: Bayesian empirical likelihood for multitask-learning: Computations and theory

Presenter: Weichang Yu, University of Melbourne, Australia

A novel Bayesian semiparametric multitask learning method is proposed. Task-relatedness is modeled using a variety of combinations of hierarchical setup and shrinkage priors to accommodate different plausible relatedness structures. To avoid misspecification of the data-generating process, Bayesian empirical likelihoods are utilized to infer both local and shared parameters. The Bayesian empirical likelihood posterior support is typically non-convex and is consequently challenging to sample from. This computational challenge is amplified in multitask settings where the likelihood arises from a combination of multiple datasets, further complicating the posterior support structure. To address this, a sequential Monte Carlo algorithm is developed that involves draws from a sequence of Bayesian posterior-adjusted empirical likelihood posterior with decreasing adjustment levels. This adjustment level sequence mitigates irregularities in intermediate steps, allowing particles to explore the posterior space efficiently. It is shown that the intermediate adjusted posterior is continuous in the tempering parameter and converges to the target original posterior as the adjustment level approaches zero, thus guaranteeing algorithmic convergence. The efficiency and convergence of the method are demonstrated through a simulation study and empirical applications.

C0633: Modelling atmospheric transport error in 4-D using the Vecchia approximation

Presenter: Michael Bertolacci, The University of Western Australia, Australia

The exchanges of CO₂ between the atmosphere and ecosystems are a vital and difficult-to-quantify component of Earth's carbon cycle. A widely used Bayesian technique for estimating these exchanges is flux inversion, where observations of atmospheric CO₂ concentrations are traced upwind to infer the fluxes. This requires a transport model that simulates the motion of CO₂ in the atmosphere. The quality of the resulting estimates is limited by the fidelity of the transport model, and current approaches to quantifying transport error are limited by computational and statistical constraints. A new method for modeling the transport error in 4-D (3-D space + 1-D time) is discussed using the Vecchia approximation. This is implemented in an existing inversion framework called WOMBAT (the Wollongong Methodology for Bayesian Assimilation of Trace-gases). Results show that this approach is practical, and simulations show that this improves the accuracy and uncertainty quantification of the flux estimates.

C0547: Revisiting parameter estimation and model selection in Gaussian mixtures of experts for clustering and regression

Presenter: Trung Tin Nguyen, Queensland University of Technology, Australia

Co-authors: Christopher Drovandi, Nhat Pham Minh Ho

Mixture of experts (MoE) models constitute a widely utilized class of ensemble learning approaches in statistics and machine learning, known for their flexibility and computational efficiency. Despite their practical success, the theoretical understanding of model selection, especially concerning the optimal number of mixture components or experts, remains limited and poses significant challenges. These challenges primarily stem from the inclusion of covariates in both the Gaussian gating functions and expert networks, which introduces intrinsic interactions governed by partial differential equations with respect to their parameters. The use of dendrograms of mixing measures is revisited, and Bayesian nonparametric techniques are incorporated to avoid predefining the number of experts. The approach enables consistent estimation of the true number of components and achieves pointwise optimal convergence rates in overfitted regimes. Importantly, it eliminates the need to train and compare multiple models with different component numbers, reducing computational costs in high-dimensional or deep learning contexts. Experiments on synthetic datasets confirm the effectiveness of the proposed method, showing superior performance over conventional criteria such as AIC, BIC, and ICL in both expert recovery and parameter estimation accuracy.

C0527: Bayesian group variable selection via penalized credible region

Presenter: Khue-Dung Dang, University of Western Australia, Australia

Co-authors: Weichang Yu

A Bayesian method is proposed for grouped variable selection in high-dimensional regression models. Most existing Bayesian methods are subjected to either high computation costs due to long MCMC runs or yield ambiguous variable selection results due to non-sparse solution output. The proposed method, GroupPenCr, is built upon the penalized credible region framework, which allows efficient computations of a sequence of sparse solutions via existing algorithms. The focus is on settings where the number of predictors grows with sample size n . For this problem, it

is proposed to use global-local shrinkage priors and to perform GroupPenCr on the mean-field variational Bayes posterior instead. This allows avoiding the computational hassle of implementing MCMC on ultra-high-dimensional data. Furthermore, it is shown that this approach is also variable selection consistent. Through extensive simulations, it is shown that GroupPenCr can outperform common methods for Bayesian group variable selection.

CO310 Room BCB 406 SCALABLE INFERENCE IN LARGE NETWORKS
Chair: Sandipan Roy
C0703: Likelihood-based methods for partially observed networks

Presenter: **Udbhav Dalavai**, Birkbeck, University of London, United Kingdom

Co-authors: Swati Chandna

Partially observed networks are common in many real-world applications due to data limitations, privacy concerns, or experimental limitations. In such settings, an appropriate statistical methodology for predicting missing links is crucial for enabling complete and meaningful analysis. The task of link prediction is studied under the random dot product graph (RDPG) model. The RDPG model provides a powerful and interpretable framework for modeling network data through low-dimensional node embeddings. It models the probability of an edge between two nodes as a dot product of their latent position vectors and is easily estimated via the well-known adjacency spectral embedding. Maximum likelihood estimation of logistic RDPGs is outlined for link prediction, community detection under partial observation, and an out-of-sample extension that makes the model suitable for time-varying networks. Empirical studies on both simulated and real-world datasets, including political blog networks and co-purchase graphs, demonstrate that the approach outperforms existing approaches.

C0753: A Bayesian approach to model uncertainty in single-cell genomic data

Presenter: **Shanshan Ren**, University College London, United Kingdom

Co-authors: Tom Bartlett, Lina Gerontogianni, Swati Chandna

Network models provide a powerful framework for analyzing single-cell count data, facilitating the characterization of cellular identities, disease mechanisms, and developmental trajectories. However, uncertainty modeling in unsupervised learning with genomic data remains insufficiently explored. Conventional clustering methods assign a singular identity to each cell, potentially obscuring transitional states during differentiation or mutation. The purpose is to introduce a variational Bayesian framework for clustering and analyzing single-cell genomic data, employing a Bayesian Gaussian mixture model to estimate the probabilistic association of cells with distinct clusters. This approach captures cellular transitions, yielding biologically coherent insights into neurogenesis and breast cancer progression. The inferred clustering probabilities enable further analyses, including differential expression analysis and pseudotime analysis. Furthermore, it is proposed to utilize the area under the curve (AUC) in clustering scRNA-seq data to quantitatively evaluate overall clustering performance. This methodological advancement enhances the resolution of single-cell data analysis, enabling a more nuanced characterization of dynamic cellular identities in development and disease.

C1050: Network estimation using graphical lasso

Presenter: **Maddie Shelley**, University of York, United Kingdom

Co-authors: Chiara Boetti, Matthew Nunes, Marina Knight

Network discovery from temporal data observed over the graph nodes has so far been investigated in both time- and spectral (Fourier)-domains, typically under the limiting assumption of weak-stationarity. The purpose is to take a spectral multiscale approach, useful in many practical contexts, under a framework that departs from stationarity. The proposed approach and its success in recovering the process dependence structure over a variety of simulated and real data scenarios are illustrated.

C1274: New directions in community detection for core-periphery networks

Presenter: **Sinyoung Park**, University of Bath, United Kingdom

Co-authors: Sandipan Roy, Matthew Nunes

Spectral clustering has been widely used as a popular tool for community detection in data with a network structure. However, spectral clustering does not perform well on certain network structures, particularly core-periphery networks. To improve clustering performance in core-periphery structures, adjacency spectral embedding (ASE) has been introduced. Despite its advantages, ASE has several limitations including its optimal performance only on dense networks. To address these limitations, a new approach is proposed, called doubled adjacency spectral embedding (DASE). It is demonstrated that DASE overcomes these challenges, particularly highlighting the improved clustering performance on both sparse and dense networks in the presence of core-periphery structures.

C1527: Performance guarantees for LLMs

Presenter: **Carey Priebe**, Johns Hopkins University, United States

Performance guarantees for LLMs are essential to myriad critical applications – enabling any application for which quantifiable confidence in LLM performance is nonnegotiable. The Data Kernel Perspective Space (DKPS) provides a framework to generate empirically validated theoretical predictions allowing guaranteed performance quantification for these transformative models. We will present the probability/statistics/linear algebra for DKPS, and illustrative examples. The probability/statistics/linear algebra for DKPS is analogous to that introduced for network time series recently.

CO210 Room BCB 408 SPORTS ANALYTICS
Chair: Andreas Groll
C0929: Parametric and semiparametric mixed modelling of athletic ability in young soccer players, considering injuries

Presenter: **Brigitte Gelein**, ENSAI, France

Co-authors: Arthur Guillotel, Benoit Bideau, Anthony Sorel

Over four consecutive seasons, a professional soccer academy was closely monitored, with data collected on key athletic performance metrics such as speed, strength, power, and endurance. To model how athletic abilities evolve with age, both parametric and semi-parametric mixed models are applied. In the semi-parametric approach, the fixed effects were estimated using machine learning methods, specifically gradient boosting and random forests. These techniques allow for greater flexibility by capturing non-linear relationships between predictors and outcomes. Importantly, the models accounted not only for age but also for injury history, a critical yet often neglected factor in athletic development. Several injury burden indicators that reflect both the duration and severity of injuries are constructed and compared in different ways. The results highlight the value of semiparametric mixed modeling approaches in producing accurate and robust predictions, with low error margins across most athletic parameters. Furthermore, the models reveal distinct developmental trajectories influenced by a combination of age and injury history. These findings provide meaningful insights to support performance monitoring, training optimization, and the long-term development of young elite soccer players.

C0937: Identifying match-fixing in professional football through in-play market dynamics

Presenter: **David Winkelmann**, Bielefeld University, Germany

Co-authors: Maya Natascha Vienken, Roland Langrock, Christian Deutscher

Match-fixing poses a significant threat to the integrity of sports by eroding public trust and undermining the commercial viability of the industry. The rise of global sports betting markets has increased opportunities for match-fixing due to the ease of access and high market liquidity, necessitating advanced detection systems. Considering the Italian Serie B, a football league historically associated with confirmed match-fixing incidents, detailed second-by-second betting volume data is utilized from a major European bookmaker to monitor market dynamics throughout a match.

The goal is to explore the potential of leveraging in-match betting data to identify suspicious behaviors among market participants that may indicate match manipulation. Specifically, time series models are applied, and residual analyses are conducted to detect deviations from typical betting patterns. The approach aims to enhance the detection of anomalies that could signify match-fixing, thereby contributing to more robust safeguarding of sports integrity.

C1182: Borrowing strength in the era of tracking data: Statistical challenges and opportunities in sport

Presenter: **Stephanie Kovalchik**, Teamworks, Germany

Borrowing strength, the use of information from related units or populations to improve estimation for a target group, is a central idea in modern statistical modeling. This principle has found particularly influential applications in sports research, where methods such as Efron and Morris empirical Bayes shrinkage and PECOTA-style projection systems stabilize player-level estimates by pooling across similar athletes. The purpose is to examine how emerging trends in data collection in both amateur and professional sport, particularly the increasing ubiquity of comparable tracking technologies and fitness monitoring tools, are generating new opportunities to extend this framework. Applications include joint modeling of training and competition data, hierarchical pooling across levels of play, and transferring information across sports with shared performance features. The major statistical and computational challenges that must be addressed are also discussed, including the need for harmonization across measurement systems, contextualization, and scalable inference methods.

C1233: Bayesian weighted discrete-time dynamic models for association football prediction

Presenter: **Roberto Macri-Demartino**, University of Trieste, Italy

Co-authors: Leonardo Egidio, Nicola Torelli

In recent years, great emphasis has been placed on the prediction of association football. Due to this, several studies have proposed different types of statistical models to predict the outcome of a football match. However, most existing approaches usually assume that the offensive and defensive abilities of teams remain static over time. The aim is to introduce a Bayesian dynamic approach for football goal-based models that uses period-specific commensurate priors to flexibly weight the evolution of attacking and defensive abilities. The approach assigns separate, time-varying precisions for each ability and period, controlled via spike and slab hyperpriors. This adaptive shrinkage borrows information about teams' strength when past and current performance aligns and allows rapid adjustments when teams experience substantial changes (e.g., transfer windows or coaching changes). This is integrated into the framework into six standard goal-based models evaluating predictive performance using data from the last five seasons of the German Bundesliga, English Premier League, and Spanish La Liga. Compared with the other discrete time dynamic models, the adaptive approach yields better predictive performance. The proposed methodology has also been implemented in the free and open source R package footBayes.

CC426 Room BCB 207 MACROECONOMIC FORECASTING AND UNCERTAINTY

Chair: Christos Savva

C1327: Survey design and professional forecasters: The case of uncertainty in the US SPF

Presenter: **Malte Kneueppel**, Deutsche Bundesbank, Germany

Co-authors: Lora Pavlova

Histogram forecasts of growth and inflation from the US Survey of Professional Forecasters (SPF) allow for an assessment of the evolution of forecast uncertainty. However, this assessment is complicated by structural breaks in measured uncertainty arising from changes in histogram bin widths. These breaks remain insufficiently addressed in the existing literature. An adjustment approach is proposed based on a structural break in 2014, during which bin widths, and consequently, measured inflation uncertainty, shifted significantly, despite true inflation uncertainty remaining virtually constant. Drawing on the results, revised bin widths to align measured uncertainty more closely with underlying uncertainty are proposed. It is recommended to differentiate bin widths between current-year and next-year forecasts.

C1387: Impacts of macroeconomic data revisions on economic activities

Presenter: **Naoko Hara**, Seikei University, Japan

Uncertainty necessitates data-driven decision-making. While sound decisions require timely and accurate economic information, major official statistics are often substantially revised long after their first release. These data revisions imply the presence of significant noise in preliminary data. The purpose is to examine the economic impacts of the noise in U.S. macroeconomic data. Noise components are identified in U.S. GDP by estimating structural vector autoregression models using different data vintages. It is found that the noise shocks significantly affect economic fundamentals and professional forecasts in the short run. Moreover, the results suggest that the GDP growth is likely to be overestimated in real time, and forecasters tend to respond differently when uncertainty increases. These findings highlight a need for caution when interpreting early releases of macroeconomic data, particularly amid heightened uncertainty.

C0198: Let the tree decide: FABART - a non-parametric factor model

Presenter: **Sofia Velasco**, European Central Bank, Queen Mary University, Germany

A novel framework integrates Bayesian additive regression trees (BART) into a factor-augmented VAR (FAVAR) to forecast macrofinancial variables and analyze oil shock asymmetries. By employing nonparametric techniques for dimension reduction, the model captures complex, nonlinear relationships between observables and latent factors that are often missed by linear approaches. A simulation experiment comparing FABART to linear alternatives and a Monte Carlo experiment demonstrate that the framework accurately recovers the relationship between latent factors and observables in the presence of nonlinearities, while remaining consistent under linear data-generating processes. The empirical application shows that FABART substantially improves forecast accuracy for industrial production relative to linear benchmarks, particularly during periods of heightened volatility and economic stress. In addition, the model reveals pronounced sign asymmetries in the transmission of oil supply news shocks to the U.S. economy, with positive shocks generating stronger and more persistent contractions in real activity and inflation than the expansions triggered by negative shocks. A similar pattern emerges at the U.S. federal state level, where negative shocks lead to modest declines in employment compared to the substantially larger contractions observed after positive shocks.

C1493: Inflation uncertainty and higher moments effects on Inflation

Presenter: **Christos Savva**, Cyprus University of Technology, Cyprus

Co-authors: Demetris Koursaros, Nektarios Michail

A new framework is developed to examine the two-way relationship between inflation and inflation uncertainty. Inflation uncertainty matters because when people are unsure about future prices, it influences wage negotiations, investment, and the credibility of monetary policy. To capture these dynamics, we adapt a flexible econometric model used initially in asset pricing to inflation. In our model, inflation depends on past inflation, inflation expectations, and the unemployment gap, while also being influenced by its own volatility and risk. Volatility is modeled using an advanced GARCH process that accounts for asymmetry and the possibility of sudden shocks. Higher-order features of the distribution, such as skewness and fat tails, are allowed to evolve over time, meaning that the model can capture both calm and turbulent periods. A key innovation is the distinction between the pure risk effect of uncertainty, which arises from volatility, and the skewness premium, which reflects how extreme shocks affect expected inflation. By combining these elements, the model shows how uncertainty not only responds to inflation but also feeds back into its expected path. This framework offers policymakers and researchers a clearer understanding of how inflation risk shapes economic outcomes and helps inform more effective monetary strategies.

CC437 Room BCB 210 FINANCIAL ECONOMETRICS I**Chair: Jose Olmo****C1168: Penalized QMLE and model selection of time series regressions****Presenter:** Julie Schnaitmann, University of Tuebingen, Germany**Co-authors:** Christian Francq, Sebastien Laurent

The purpose is to examine a linear regression model applied to the components of a time series, aiming to identify time-varying, constant as well as zero conditional beta coefficients. To address the non-identifiability of parameters when a conditional beta is constant, a Lasso-type estimator is employed. This penalized estimator simplifies the model by shrinking the estimates in favor of natural constant beta representations. A multistep estimator that first captures the dynamics of the regressors is proposed before estimating the dynamics of the betas. This strategy breaks down a large-dimensional optimization problem into several lower-dimensional ones. Since making strict parametric assumptions is avoided about the innovation distributions, quasi-maximum likelihood estimators are used. The non-Markovian nature of the global model means that standard convex optimization results cannot be applied. The asymptotic distribution of the multistep Lasso estimator and its adaptive version are analyzed, deriving bounds on the maximum value of the penalty term. A nonlinear coordinate-wise descent algorithm is also proposed, which is demonstrated to find stationary points of the objective function. The finite-sample properties of these estimators are further explored through a Monte Carlo simulation and illustrated with an application to financial data.

C1206: Expected shortfall regression for high-dimensional additive models**Presenter:** Toshio Honda, Hitotsubashi University, Japan**Co-authors:** PoHsiang Peng

The expected shortfall (ES) regression is a useful tool to analyze the relation between the response variable and the covariates through quantile and conditional mean. As is well-known, there is no single loss function for expected shortfall estimation. Recently, a two-step procedure for ES regression was proposed, and this is successful due to the Neyman orthogonality. Then, based on the findings, high-dimensional linear ES regression models and nonparametric ES models were considered. To tackle both non-linearity and high-dimensionality, additive models are assumed for both quantile and expected shortfall in the high-dimensional settings, and the group Lasso and SCAD estimators are considered. The oracle inequality and the oracle property are established for them. The theoretical results imply that quantile estimation does not affect ES estimation asymptotically. Numerical results that demonstrate satisfactory performances in model selection, estimation accuracy, and prediction error are presented for a moderate sample size.

C1320: Synthetic control method with mixed frequency data**Presenter:** Lu Zhang, University of Science and Technology of China, China

Mixed-frequency data, where variables are observed at different temporal resolutions, commonly occur in economic and financial studies. Classical synthetic control methods (SCM) are ill-suited for such data, often necessitating aggregation or prefiltering that may discard valuable information. The aim is to propose a novel mixed-frequency synthetic control method (MF-SCM) to integrate mixed-frequency data into the synthetic control framework effectively. A flexible estimation procedure is developed to construct synthetic control weights under mixed-frequency settings and establish the theoretical properties of the MF-SCM estimator. Specifically, it is first proven that the estimator achieves asymptotic optimality, in the sense that it achieves the lowest possible squared prediction error among all potential treatment effect estimators from averaging outcomes of control units. The asymptotic distribution of the average treatment effect (ATE) estimator is then derived, using projection theory, and confidence intervals are constructed for the ATE estimator. The method's effectiveness is demonstrated through numerical simulations and two empirical applications on air pollution alerts and a policy study on the Tax Cuts and Jobs Act of 2017 in the US.

C1252: Bootstrap-based inference for weak instruments in IV regression in finance**Presenter:** Yongdeng Xu, Cardiff University, United Kingdom

Econometric work on weak instruments typically analyzes a worst-case endogeneity setting ($\rho = 1$) with symmetric two-sided tests, under which 2SLS t-tests over-reject. In empirical finance, endogeneity is usually small to moderate; standard t-tests are instead conservative and, crucially, their sampling distribution is skewed, with rejections concentrated in one tail. Motivated by this gap, restricted-bootstrap t-tests that impose the null during resampling are developed. In extensive Monte Carlo experiments, a restricted wild efficient bootstrap attains nominal size and correct tail behavior across a wide range of instrument strengths and endogeneity levels, addressing both conservativeness and asymmetry. The approach also extends the ViF idea of a prior study beyond the just-identified case by accommodating overidentified models with multiple instruments and delivering valid one- and two-sided inference. For practice, unified critical-value tables indexed only by the first-stage F and the endogeneity level (ρ) are provided, together with MATLAB/Stata code. These tools are drop-in replacements for standard t-tests, enabling reliable 2SLS inference.

CC360 Room BCB M201 ROBUST METHODS**Chair: Maria Brigida Ferraro****C0190: Robust beta regression through the logit transformation****Presenter:** Francisco F Queiroz, University of Sao Paulo, Brazil**Co-authors:** Silvia Ferrari, Yuri Maluf

Beta regression models are employed to model continuous response variables in the unit interval, like rates, percentages, or proportions. Their applications rise in several areas, such as medicine, environmental research, finance, and natural sciences. The maximum likelihood estimation is widely used to make inferences about the parameters. Nonetheless, it is well-known that the maximum likelihood-based inference suffers from the lack of robustness in the presence of outliers. Such a case can bring severe bias and misleading conclusions. Recently, robust estimators for beta regression models were presented in the literature. However, these estimators require non-trivial restrictions in the parameter space, which limit their application. New robust estimators that overcome this drawback are developed. Their asymptotic and robustness properties are studied, and robust Wald-type tests are introduced. Simulation results evidence the merits of the new robust estimators. Inference and diagnostics using the new estimators are illustrated in an application to health insurance coverage data.

C0282: Novel robust estimator in logistic regression**Presenter:** Alfonso Garcia-Perez, Universidad Nacional de Educacion a Distancia (UNED), Spain

Down-weighting the observations according to their outlyingness, understood as a large distance to the center of the data, is an old idea in Robustness. In fact, the famous robust regression estimator LTS (trimmed least squares) and, in general, all trimmed estimators weight with 0 or 1 the ordered observations. There are some problems with this approach, mainly related to the difficulty in multivariate statistics to order the observations. The aim is to propose to change the concept of outlyingness, understood here as the low probability of obtaining the value of the estimator in the problem considered. These ideas are applied to define the novel distribution-weighted estimator in logistic regression. Its main properties are also studied.

C1172: Extending TCLUSST to higher dimensions**Presenter:** Luis Angel Garcia-Escudero, Universidad de Valladolid, Spain**Co-authors:** Lucia Trapote-Reglero, Agustin Mayo-Isacar

Outliers are known to significantly distort the results of many commonly used clustering methods, often leading to unreliable cluster partitions. To address this issue, different robust clustering approaches have been developed that not only reduce the influence of but also facilitate the detection

of meaningful outliers. The focus is on robust clustering methods based on trimming, especially TCLUST, which extends the type of trimming used by MCD in one-population problems to allow for different subpopulations or clusters unknown in advance. While TCLUST performs well on low-dimensional data, it struggles with high-dimensional datasets due to the complexity involved in estimating a large number of parameters. The robust linear grouping (RLG) method offers an alternative by assuming clusters lie near lower-dimensional subspaces, thus combining clustering with dimensionality reduction. However, RLG has limitations when subspaces intersect and assumes simplistic isotropic orthogonal errors. A robust clustering method extending TCLUST is presented, which builds on the high-dimensional data clustering (HDDC) method by including trimming and eigenvalue constraints. This approach balances TCLUST and RLG, requiring careful adaptation of TCLUST and HDDC steps for proper implementation. An extension allowing for cellwise trimming is also outlined.

C1486: Byzantine-robust distributed one-step estimation

Presenter: **Chuhan Wang**, Beijing Normal University, China

A Robust One-Step Estimator (ROSE) is proposed that solves the Byzantine failure problem in distributed M-estimation with only a single iteration, even when a moderate fraction of worker nodes exhibit Byzantine behavior. This estimator achieves higher asymptotic relative efficiency than conventional median-based estimators, while maintaining the same computational complexity. It is also robust to outliers or missing data that may occur during centralized aggregation. We establish the asymptotic normality of the estimator as the parameter dimension p increases with the sample size, and under mild assumptions, we derive its convergence rate.

CC430 Room BCB 407 STATISTICAL ALGORITHMS AND MACHINE LEARNING TECHNIQUES

Chair: Jonas Andersson

C1381: Interpretnn: An R package for statistically-based neural network interpretation

Presenter: **Andrew McInerney**, University of Limerick, Ireland

Co-authors: Kevin Burke

Feedforward neural networks (FNNs) are typically viewed as pure prediction algorithms, and their strong predictive performance has led to their use in many machine-learning applications. Their success in predictivity can be attributed, at least in part, to their ability to capture complex relationships through the modelling of higher-order interactions. However, their flexibility comes with an interpretability trade-off; thus, FNNs have been historically less popular among statisticians, who tend to use more interpretable additive models. Nevertheless, classical statistical theory, such as significance testing and uncertainty quantification, is still relevant for FNNs. Supplementing FNNs with methods of statistical inference, model selection, and covariate-effect visualizations can shift the focus away from black-box prediction and make FNNs more akin to traditional statistical models. This can pave the way towards more inferential analyses. The focus is on the use of the R package, interpretnn, which extends existing neural network objects in R to allow for more statistically-based outputs. The aim of this package is to improve interpretation and to increase the utility of the neural network for statisticians.

C1465: Modelling competing risks data through clustering of cumulative incidence functions

Presenter: **Marta Sestelo**, University of Vigo, Spain

Co-authors: Nora M Villanueva, Luis Machado, Javier Roca Pardinas

The cumulative incidence function is the standard tool for estimating the marginal probability of a given event in the presence of competing risks. One basic but important goal in the analysis of competing risk data is the comparison of these curves, for which limited literature exists. An R package is presented that implements a procedure to not only test the equality of cumulative incidence curves but also cluster them when differences exist. The package automatically determines group composition and selects the optimal number of clusters. The applicability of the proposed method is illustrated using real data. This tool provides researchers with a practical and accessible framework for exploring and analyzing competing risks data.

C1419: Algorithmic modeling of a complex k-out-of-n system with corrective and preventive repair and Bernoulli vacation policy

Presenter: **Mohammad Dawabsha**, ARAB AMERICAN UNIVERSITY, Palestine

Co-authors: Juan Eloy Ruiz-Castro

The algorithmic modeling and optimization of complex reliability systems is of great interest in this field. In this work, a k-out-of-n: G system is modeled in which multi-state units are subject to multiple events, preventive maintenance, and a vacation policy of the repair facility. Each unit may experience either repairable or non-repairable internal failures, as well as external shocks. Each external shock to the system can cause wear in the units, an extreme failure, or a modification in their internal operating behavior. Preventive maintenance is also introduced when a critical level of degradation is reached. All events are communicated to the repairer through monitoring. Therefore, there are two tasks to be carried out in the repair channel: corrective repair and preventive maintenance. To optimize the system, a Bernoulli vacation policy is introduced in the repair facility, depending on the number of operational units within it. The system is algorithmically modeled using matrix-analytic methods and Markovian arrival processes with marked arrivals. Multiple performance measures are constructed, and costs and associated measures are introduced. The entire study is conducted under both transient and stationary regimes using matrix-analytic methods. For model optimization, a Pareto multi-objective optimization framework is introduced. An example is presented to show the versatility of the work.

C1391: Systematic review: Comparison between generalized additive models and neural networks across application areas

Presenter: **Jessica Doohan**, University of Limerick, Ireland

Co-authors: Kevin Burke

Over the past two decades, advancements in machine learning, particularly neural networks, have reshaped predictive modelling. While neural networks are often regarded as black-box algorithms, their universal approximation properties have made them highly effective predictors, albeit difficult to interpret. Much of this development has occurred outside statistics, creating a disconnect in terminology and methodology. Recent works have highlighted parallels with traditional statistical approaches, showing that multilayer perceptrons can be viewed as nonparametric generalisations of regression models. Despite these connections, comparative studies have largely focused on simple models such as linear regression, often within specific application domains. Such comparisons are yet to consider more flexible statistical models. Generalized additive models (GAMs), which extend linear regression by allowing non-linear predictor effects while preserving interpretability, provide a more appropriate benchmark. To date, no systematic review has examined how GAMs perform relative to neural networks. The focus is on that gap through a systematic review of over 140 papers and 400 datasets. Beyond summarizing comparisons, reported performance metrics are analyzed using mixed-effects modelling to investigate characteristics that can explain and quantify observed differences, including application area, publication year, sample size, number of predictors, and neural network complexity.

CC419 Room BCB 409 SPATIAL STATISTICS

Chair: Enea Bongiorno

C1023: Down to earth: Estimating GDP components in small areas using satellite data

Presenter: **Luca Romagnoli**, University of Molise, Italy

Co-authors: Claudio Lupi

Although the development of macroeconomic applications of data extracted from satellite imagery is relatively recent, satellite data have been regularly used to derive GDP estimates, especially in regions where official statistics are scarce or unavailable. Satellite nighttime light and land cover data are used here together with cross-validated geographically weighted regressions to estimate disposable income at the municipal level in Italy, two years ahead of official statistics. The geographically weighted regressions are used explicitly to account for the known spatial non-

stationarity of the relationship between nighttime light and economic activity. In addition to disposable income at the municipal level, an estimate of value added with the same spatial breakdown is produced to supplement the official statistics, which are limited to the level of the Italian provinces and are only available with a considerable delay. For larger municipalities, finer estimates are possible, which map the spatial heterogeneity of income and value added even within the same municipality.

C1246: Detecting disease clusters using functional additive models with spatial dependence

Presenter: **Michio Yamamoto**, The University of Osaka / RIKEN AIP / Shiga University, Japan

Co-authors: Tatsuhiko Anzai, Kunihiko Takahashi

Identifying spatial disease clusters is essential for characterizing disease patterns and informing prevention and treatment strategies. Spatial scan statistics are widely used for cluster detection with a variable scanning window size. When covariates influence the outcome and are not randomly distributed across space, cluster searches should adjust for them. Moreover, spatial correlation in the outcome, which is often overlooked in practice, can affect detection. The aim is to propose a new spatial scan statistic that accommodates multiple functional covariates summarizing longitudinal histories and explicitly accounts for spatial correlation in the outcome. These factors are modeled flexibly within a functional additive modeling framework. An optimization algorithm is developed for parameter estimation under Gaussian outcomes. Simulation studies and an application to real data show that, compared with existing methods, the proposed approach reliably detects disease clusters in the presence of longitudinal covariates and spatial correlation.

C1263: Spatial detection of adjacent hotspot clusters in disease mapping

Presenter: **Kunihiko Takahashi**, Institute of Science Tokyo, Japan

Co-authors: Hideyasu Shimadzu

In spatial epidemiology, statistical tests are often employed to detect regional disease clusters, particularly to evaluate whether disease risk is significantly elevated compared to surrounding areas. A central method is the cluster detection test (CDT), which identifies non-random spatial distributions and highlights high-prevalence regions without prior assumptions. Among CDT approaches, scan statistics based on maximum likelihood ratios, such as Kulldorff's circular scan and Tango and Takahashi's flexibly shaped scan, have been widely applied. More recent developments enable the simultaneous detection of multiple clusters by integrating generalized linear models with information criteria. However, conventional scan-based methods often assume uniform risk within a single cluster, leading to the erroneous merging of adjacent hotspots with different risk levels. To address this limitation, a new scan-based procedure is proposed, incorporating Cochran's Q-statistic to evaluate heterogeneity within clusters. This extended framework offers the accurate identification of adjacent hotspots as distinct clusters. Through real-world applications, the improved performance of the proposed method compared with existing scan-based tests is demonstrated.

C1337: Vecchia-inducing-points full-scale approximations for Gaussian processes

Presenter: **Tim Gyger**, Lucerne University of Applied Sciences, Switzerland

Co-authors: Fabio Sigris, Reinhard Furrer

Gaussian processes are flexible, probabilistic, non-parametric models widely used in machine learning and statistics. However, their scalability to large data sets is limited by computational constraints. To overcome these challenges, Vecchia-inducing-points full-scale (VIF) approximations are proposed, combining the strengths of global inducing points and local Vecchia approximations. Vecchia approximations excel in settings with low-dimensional inputs and moderately smooth covariance functions, while inducing point methods are better suited to high-dimensional inputs and smoother covariance functions. The VIF approach bridges these two regimes by using an efficient correlation-based neighbor-finding strategy for the Vecchia approximation of the residual process, implemented via a modified cover tree algorithm. The framework is further extended to non-Gaussian likelihoods by introducing iterative methods that substantially reduce computational costs for training and prediction by several orders of magnitude compared to Cholesky-based computations when using a Laplace approximation. In particular, novel preconditioners are proposed and compared, and theoretical convergence results are provided. Extensive numerical experiments on simulated and real-world data sets show that VIF approximations are both computationally efficient, as well as more accurate and numerically stable than state-of-the-art alternatives.

CP001 Room BCB 208 POSTER SESSION

Chair: Louisa Kontoghiorghe

C0234: Sparse approximation of kernel function estimators via genetic algorithm

Presenter: **Kihei Nishida**, Kyoto Sangyo University, Japan

A method is proposed for constructing multivariate kernel-based function estimators, specifically, both kernel density and regression estimators, using a genetic algorithm to obtain sparse representations. The method applies sparsification in two forms: Reducing the number of data points (data size sparsification) and reducing the number of input dimensions (feature sparsification). The algorithm generates multiple subsamples of user-specified fixed size through random sampling with replacement. Each subsample is treated as a chromosome, and each gene within it, in the terminology of genetic algorithms, represents either a data point or a variable, depending on the selected sparsification mode. Pairs of subsamples undergo genetic operations such as crossover, mutation, and direct inheritance, applied with predetermined probabilities. Fitness is evaluated based on performance in approximating the target density or regression function, and those with superior fitness are selected to survive and contribute to the next generation. Through repeated iterations, the algorithm evolves toward a compact yet accurate estimator. Simulation results show that the proposed method consistently outperforms well-known kernel-based estimators in accuracy and data condensation ratio.

C0345: A novel statistical framework for case-control genome-wide association studies

Presenter: **Itziar Irigoien**, University Basque Country, Spain

Co-authors: Selena Aranda, Marina Mitjans, Bru Cormand, Concepcion Arenas, Group Cibersam

Genome-wide association studies (GWAS) investigate the relationship between genetic variation and traits of interest by analyzing molecular markers distributed throughout the genome, such as single-nucleotide polymorphisms (SNPs). The use of GWAS has been growing over the years, enabling the identification of genetic loci associated with a wide range of phenotypes, including asthma, diabetes, and psychiatric disorders. In case-control designs, the most commonly used statistical method is logistic regression, which typically includes as covariates the principal components derived from a multidimensional scaling (MDS) analysis to account for population stratification. The optimal number of components to include depends on the population structure and sample size, although the inclusion of up to 10 components is generally accepted. Given that many SNPs are correlated due to linkage disequilibrium, the Bonferroni correction for multiple testing is often overly conservative, leading to an increased risk of false negative findings. To address this limitation, a multivariate statistical approach is proposed based on distances and k-neighborhoods, which offers advantages in detecting associations in case-control GWAS. The application of this method is illustrated using real data, highlighting both its implementation and the results it provides.

C1128: Estimating all-cause excess mortality in the Philippines (2020 to 2023) using the Farrington algorithm

Presenter: **Rutcher Lacaza**, University of the Philippines, Philippines

Excess mortality estimates have been widely used to assess the impact of COVID-19 on deaths during and after the pandemic. Excess deaths in the Philippines from 2020 to 2023 were estimated using weekly all-cause deaths data from the Philippine Statistics Authority, with a focus on the later stages of the pandemic and the post-emergency period. The Farrington model, a quasi-Poisson regression approach, was applied to account for seasonal variation and overdispersion. Excess deaths were assessed across two periods: During the pandemic (1 January 2020 to 7 May 2023) and after the lifting of the emergency (8 May to 31 December 2023), disaggregated by age, sex, and province. Results show no significant excess mortality in 2020, followed by a rise in 2021, and peaking in 2022 before declining in 2023. Excess deaths were highest among those aged 60 and

above, with males slightly more affected than females. By 2022, all provinces recorded positive excess mortality. Comparing the post-emergency period with prior years shows the highest percentage of excess mortality in 2022. The findings highlight the importance of continuous mortality monitoring, refinement of statistical methods, and improvements in death registration systems to fully assess the extent of epidemics on mortality.

C1265: Tobit filtering for nonlinear systems under random access protocol: Dealing with multirate measurements

Presenter: **Antonia Oya**, Universidad de Jaen, Spain

Co-authors: Shuo Yang, Jun Hu, Raquel Caballero-Aguila

Over the past years, growing research interest in filtering for networked systems has been driven by a number of crucial challenges. These include limited communication resources, heterogeneous sampling rates, and imperfect measurements. In many practical scenarios, sensor data may be censored as a direct result of resolution constraints or energy-saving strategies. This issue can be further complicated in large-scale systems, where a variety of heterogeneous devices naturally lead to the coexistence of multiple sampling periods. Moreover, when numerous sensors attempt to transmit data simultaneously, random access protocols are often adopted to alleviate communication congestion and reduce collisions. In this context, a Tobit filtering approach is proposed for multirate nonlinear systems with censored measurements operating under a random access protocol. More specifically, an innovative filtering scheme is developed, where a minimized upper bound of the filtering error covariance matrix is obtained, and a recursive procedure is established to determine the filter gain matrix. Furthermore, a sufficient condition for the boundedness of the filtering error is provided from an analytic perspective. Finally, some experimental comparisons are conducted to validate the effectiveness and highlight the advantages of the proposed Tobit filtering algorithm.

C1272: A cross-validation bandwidth choice for nonparametric tests in regression models

Presenter: **Sandie Ferrigno**, INRIA, France

Co-authors: Marie-Jose Martinez

Many goodness-of-fit tests assess the different assumptions of a regression model. The focus is on the task of choosing the structural part of the variance function in the regression model, and three nonparametric tests are considered, all based on generalizations of the Cramer-von Mises statistic. Nonparametric kernel methods are used for regression and variance function estimations. Bandwidth selection is of key practical importance for these estimations. The bandwidths for the regression and variance functions are obtained separately by cross-validation methods. A simulation study, based on wild bootstrap methods, is carried out to compare the three tests in terms of statistical significance and power function.

C1297: Prognostic diagnostics of positivity violation in causal inference for longitudinal multilevel data

Presenter: **Huixia Savannah Wang**, Umea University, Sweden

Positivity is a fundamental assumption in causal inference, requiring that all individuals have a nonzero probability of receiving each exposure level given their covariates. While the propensity score is commonly used to assess this assumption through balancing covariates with respect to exposure, the prognostic score instead summarizes covariate associations with potential outcomes. Prognostic score adjustment reduces outcome variation between groups, leaving only variation due to the exposure effect and random noise. Although the prognostic score is naturally aligned with the estimation of the average exposure effect under no exposed group, its use is extended as a diagnostic and trimming tool to support the estimation of the average exposure effect in a longitudinal setting. The approach applies prognostic score adjustment to ensure comparability, then fits a random-intercept Bayesian additive regression tree (riBART) model for flexible outcome modeling, and finally estimates the average exposure effect via g-computation. This implementation helps mitigate challenges arising from the positivity violation. The approach is further applied to SHARE data, examining the causal effect of depression on human cognitive function.

C1377: A general model for partially observed functional data

Presenter: **M Carmen Aguilera-Morillo**, Universitat Politècnica de Valencia, Spain

Co-authors: Maria Durban, Pavel Hernandez

In functional data analysis, it is usual to assume that all sample data are fully observed within their domain. However, there are situations where the sample data (curves, surfaces) are partially observed, i.e., contain gaps or missing parts. To deal with the problem of statistical modeling of partially observed multidimensional functional data, a generalized additive scalar-on-function regression model is proposed. In order to control the smoothness of the functional coefficient, a p-spline penalty is added to the model estimation. The functional model is estimated thanks to the connection with the mixed effects model, where the smoothing parameters and the model coefficients are estimated directly. The performance of the proposed model is tested on a simulated and a real dataset of images of air pollution from India.

C1380: On survival trees for interval-censored survival data

Presenter: **Asanao Shimokawa**, Tokyo University of Science, Japan

Generally, survival time analysis often considers only right-censored data. In the construction of survival trees, many studies also assume right-censoring. However, in practice, many data sets are more appropriately modeled under the assumption of interval-censoring rather than right-censoring. Accordingly, the focus is on the construction of survival trees that account for interval censoring. For this purpose, approaches are proposed based on hypothesis tests as well as on prediction accuracy metrics. The performance of this criterion is evaluated through simulation studies. In addition to this, the results of the survival tree obtained from actual data are shown.

C1423: Global high frequency price forecasting through distributed consensus using local microstructure

Presenter: **Juan Francisco Munoz-Elgueazabal**, Western Institute of Technology and Higher Education (ITESO), Mexico

Co-authors: Juan Diego Sanchez-Torres

The purpose is to introduce a novel distributed convex optimization framework for high-frequency financial forecasting that addresses the fundamental econometric challenge of achieving global consensus across heterogeneous market environments while preserving local adaptation capabilities. The methodology extends classical diffusion strategies through dynamic regularization mechanisms, implementing both combine-then-adapt (CTA) and adapt-then-combine (ATC) protocols within leader-follower network topologies. The framework provides theoretical convergence guarantees under convexity assumptions while maintaining computational tractability for millisecond-level market microstructure analysis. Empirical validation using high-frequency cryptocurrency market data across nine distributed agents processing 77.76 million observations demonstrates superior econometric performance: 21% accuracy improvement over traditional centralized approaches with sub-100 millisecond processing latency. The distributed architecture exhibits remarkable robustness, showing less than 8% performance degradation during regional node failures, making it particularly suitable for fragmented cryptocurrency markets where traditional econometric models fail to capture cross-exchange dependencies and regional microstructure heterogeneity effectively.

C1444: Studying urban well-being through principal component analysis

Presenter: **Olegs Krasnopjorovs**, University of Latvia, Latvia

In a world displaying a plethora of information, composite social and economic indicators have become extremely popular as a way to aggregate different dimensions in a single index, dimensions that would otherwise be difficult or impossible to compare. Eurobarometer survey on the quality of life is a rich dataset containing responses for about 40 urban well-being measures collected in 83 European cities from over 70 thousand respondents. The purpose is to show how the application of principal component analysis (PCA) to the Eurobarometer survey data allows dimensionality decrease of this large dataset without a notable information loss. PCA is employed to obtain one numerical variable for each area of urban well-being (safety, trust, environment, infrastructure, public transport, governance, livability, and economic situation) and a composite numerical urban well-being indicator per city (based on the value of first principal component of the PCA), which are then rescaled to a 0-100 point scale

for representation purposes. The use of PCA is justified by the Kaiser-Meyer-Olkin criterion (which has a value of 0.85 in the composite urban well-being indicator). Furthermore, the relationship between the composite urban well-being indicator and various city characteristics is explored, and its positive relationship with gross domestic product per capita in a city as well as various city amenities, as well as its negative relationship with city size and congestion are shown.

C1498: Statistical reproducibility for Tukey's honest significant difference test

Presenter: **Norah Alshahrani**, University of Bisha, Saudi Arabia

Reproducibility is a cornerstone of reliable research and has received increasing attention in recent years. Statistical reproducibility, in particular, concerns whether statistical conclusions remain consistent when experiments are replicated under identical conditions. This paper investigates the reproducibility probability (RP) of Tukey's Honest Significant Difference (HSD) test, a widely used post hoc procedure following one-way ANOVA, within the framework of nonparametric predictive inference (NPI). Simulation results show that RP declines as p-values approach the significance threshold, indicating instability in such cases. Moreover, sample size exerts a notable influence: larger samples decrease RP when $p\text{-value} > \alpha$, but increase RP when the $p\text{-value} < \alpha$. Distributional characteristics also affect RP, especially when groups have unequal variances. These findings highlight the importance of jointly considering p-values, sample size, and distributional features when interpreting the reproducibility of post hoc test results.

C1519: Optical character recognition enhancement for manufacturing instrument panel numbers with a new object detection loss

Presenter: **SeokHwan Hong**, Inha University, Korea, South

Co-authors: Donghyeon Yu

Optical character recognition is widely used in various fields, such as license plate recognition and text extraction from images. In particular, number recognition performs well when the target area is known and fixed. However, unlike license plate recognition, recognizing manufacturing instrument panel numbers remains challenging when the recognition areas vary in location and size. We propose a new loss function for object detection to enhance optical character recognition of manufacturing instrument panel numbers. The proposed loss function enables the object detection method to identify the smallest recognition area that contains all manufacturing instrument panel numbers. Using real images obtained from manufacturing factories, we show that the proposed object detection loss function improves the optical character recognition model for instrument panel numbers.

C1520: Enhancing nighttime road image segmentation with CycleGAN-based image transformation

Presenter: **Jiho Lee**, South Korea, Inha University, Korea, South

Co-authors: Donghyeon Yu

Image segmentation plays an important role in image recognition in computer vision, as it provides not only object labels but also the boundaries of the objects. However, annotated datasets are insufficient compared to other image recognition tasks, such as object detection with bounding boxes, since datasets for image segmentation require pixel-level object labels. In particular, it is more difficult to annotate pixel-level object labels for images obtained at night. In addition, the performance of image recognition decreases for nighttime images, which are affected by noise from car headlights and streetlamps. In this study, we propose a novel framework that integrates cycleGAN-based image transformation with state-of-the-art segmentation models. Experiments on road images demonstrate that our approach significantly enhances segmentation performance for nighttime images.

C1511: Truncated inverse-Levy measure representation of the beta process

Presenter: **Junyi Zhang**, The Education University of Hong Kong, Hong Kong

Co-authors: Angelos Dassios, Chong Zhong, Qiufei Yao

The beta process is a widely used nonparametric prior in Bayesian machine learning. While various inference schemes have been developed for the beta process and related models, the current state-of-the-art method relies heavily on the stick-breaking representation with decreasing atom weights, which is available only for a special hyperparameter. In this work, we introduce the truncated inverse-Levy measure representation (TILe-Rep) that extends the decreasing atom weights representation of the beta process to general hyperparameters. The TILe-Rep fills the gap between the two previous stick-breaking representations. Moreover, it has lower truncation error than other sequential representations of the beta process and may lead to the posterior consistency property of Bayesian factor models. We demonstrate the usage of the TILe-Rep in the celebrated binary latent feature model and the beta process factor analysis model.

Sunday 14.12.2025

10:55 - 12:10

Parallel Session G – CFE-CMStatistics 2025

CO028 Room BCB G07 HiTEC: STATISTICAL INFERENCE ON INSTABILITIES AND INDEPENDENCE**Chair: Matus Maciak****C1155: Online change-point detection: A traders perspective***Presenter:* **Bojana Milosevic**, University of Belgrade, Serbia*Co-authors:* Zikica Lukic

The aim is to present several classes of integral-transform-based test statistics for change-point analysis, examining their asymptotic behavior and small-sample performance. These methods are then adapted for real-time monitoring, enabling the rapid detection of structural changes as new data arrive. From a trader's perspective, such tools are critical for reacting to sudden regime shifts. Practical implementation details are highlighted, and it is illustrated how these methods can support decision-making in unstable financial environments.

C1212: Online detection of risk instabilities based on conditional expectiles*Presenter:* **Matus Maciak**, Charles University, Czech Republic*Co-authors:* Michal Pesta, Gabriela Ciuperca

The purpose is to consider an online changepoint detection procedure based on conditional expectiles, which are convenient tools for risk assessment tasks in empirical finance. The expectiles are well-known in econometrics for being a coherent and elicitable risk measure. In addition, the approach based on the expectiles introduces some robustness when compared with traditional moment-based techniques, and it also provides a more complex insight into the overall data-generating mechanism. The proposed statistical test is proven to be consistent, while the distribution under the null hypothesis does not depend on the underlying functional form or the unknown parameters. Theoretical details and finite sample performance are discussed in the talk.

C1235: On independence testing in the presence of cure data*Presenter:* **Marija Cuparic**, University of Belgrade, Serbia*Co-authors:* Bojana Milosevic

The focus is on the problem of testing independence between two randomly right-censored variables in the presence of cure data, modeled through an additive mixture framework with a cured fraction. Several classes of test statistics are proposed, inspired by Kendall's tau and its extensions. Their asymptotic properties are derived, and their finite-sample performance is thoroughly investigated. Robustness to different censoring models is also explored. In addition, a potential generalization that can be used in more dimensional settings is presented. A comprehensive simulation study demonstrates the strong competitiveness of the proposed procedures.

C1500: Computational methods and algorithms based on characteristic functions and their implementation in CharFunTool*Presenter:* **Viktor Witkovsky**, Slovak Academy of Sciences, Slovakia

Characteristic-function-based (CF) methods provide a powerful analytical framework for statistical inference, yet their numerical implementation is often affected by oscillatory integrals, truncation, and loss of precision. This work presents stabilized computational techniques implemented in the open MATLAB framework CharFunTool, improving robustness and accuracy of CF evaluation and inversion. The approach combines Fourier inversion with exponential damping and smooth windowing, double-exponential (DE) quadrature for rapidly convergent integration over infinite domains, and adaptive quadrature guided by local oscillation analysis. Regularized and weighted estimators of empirical CFs further suppress high-frequency noise. The results demonstrate that numerically stabilized CF-based algorithms provide a coherent and robust framework for inference under possible instability and analytical intractability.

CO188 Room BCB G08 PHYSICS-GUIDED SPATIOTEMPORAL DATA: DESCRIPTIVE VS DYNAMICAL SYSTEMS **Chair: Arthur Guillaumin****C1189: Flexible hierarchical Bayesian process-informed neural network models for spatiotemporal data***Presenter:* **Christopher Wikle**, University of Missouri, United States

Process (physics)-informed neural models have become ubiquitous across many areas of science in recent years due to the value of regularizing neural networks based on an underlying partial differential equation physical constraint. The purpose is to discuss a generalization in which mechanistic information about the process can be incorporated via a Bayesian hierarchical approach. In many scientific applications where there is substantial a priori process knowledge, incorporating this information can improve model performance and efficiency. The notion of including process knowledge in data-driven models for spatiotemporal data is not new (e.g., data assimilation, physical-statistical modeling, etc.), and considering hybrid statistical/neural approaches can provide more realistic modeling of complex processes while quantifying uncertainty. A brief overview of neural and statistical approaches is presented, and a unifying hierarchical modeling structure is presented that can accommodate flexible, mechanistically informed neural or statistical models for spatiotemporal dynamic processes.

C1203: On space-time local interaction machines*Presenter:* **Dionissios Hristopulos**, Technical University of Crete, Greece

The focus is on the development of local interaction machines (SLIMs) for scalar and vector (multi-output) spatiotemporal processes. SLIMs are inspired by statistical physics models and employ sparse precision matrices. Hence, they can significantly improve the computational efficiency of regression tasks compared to the Gaussian process framework. The construction of positive-definite precision matrices that can incorporate general distance measures is discussed. Parameter estimation and regression equations based on these precision matrices are presented, and example datasets are used to illustrate the application of SLIMs.

C1467: Physics-guided space-time mapping of sea surface height: Novel covariance models and methods*Presenter:* **Arthur Guillaumin**, Queen Mary University of London, United Kingdom*Co-authors:* Cimarron Wortham, Early Jeffrey, Cheuk Yan Cliff Lo

The purpose is to design a space-time covariance model of sea surface height of the oceans by leveraging large gridded simulation data from an idealized ocean model run in various physical regimes. Two approaches are considered: one parametric, which relies on a mixture of Lagrangian space-time covariance models, and another one akin to a non-parametric approach, where a deep neural network is trained to directly infer the frequency-wavenumber spectrum of the process. Some novel approaches are also discussed for efficient approximate inference, both for the estimation of parametric covariance models and for Gaussian prediction. The models and methods are applied to the mapping of sea surface height measurements: Firstly, on simulation data from a coupled climate model, and secondly, to real satellite data.

CO069 Room BCB G09 FINANCIAL AND MACRO ECONOMETRIC PREDICTABILITY**Chair: Robinson Kruse-Becher****C0678: Signal-based subset selection for long-term portfolio returns***Presenter:* **Rainer Alexander Schuessler**, University of Muenster, Germany*Co-authors:* Sven Lehmann, Mark Trede

A simple framework is proposed where sparse equal-weighted subsets enhance long-term portfolio performance when return signals are weak and noisy. Building on a prior study, which shows that even randomly selected, equal-weighted small portfolios can outperform the market when rebalanced, this insight is translated into a systematic signal-based selection that amplifies the effect. The design combines simple predictive

signals, such as short-term momentum with L0-based subset selection and equal weighting, to contain estimation error and shift the bias-variance trade-off in favor of more robust realized long-term compound returns.

C0692: **Macroeconomic survey forecasting in times of crises**

Presenter: **Robinson Kruse-Becher**, FernUniversität in Hagen, Germany

Co-authors: Philip Letixerant

Survey-based forecasts such as the SPF are generally of great importance. The aim is to investigate the potential for improvement by exploiting historical information that is similar to the current situation in which the expert forecasts are made. The real-time adjustment applied to the SPF is a similarity-based intercept correction. Periods of similarity are found by a nearest neighbor matching procedure. A large number of aspects in the matching algorithm are investigated by considering, e.g., matching variables (levels versus forecast errors), matching quality measures, and weighting based on the degree of similarity, a recency filter, etc. Hyperparameters such as the block length and the number of matched periods to be averaged are selected by cross-validation. The focus is on the great financial crisis and Covid-19 recession, and the potential to reduce the relative out-of-sample MSFE is analyzed in such times of crises. The analysis is complemented by considering the performance of the adjustment in non-crisis periods.

C0778: **Managerial habit formation and predictability**

Presenter: **Christoph Wegener**, Leuphana University Lüneburg, Germany

Within a framework featuring managerial habit formation, closed-form expressions are derived for predictive regressions, in which the dividend-price ratio forecasts future dividend growth rates and returns. It is found that return predictability can arise from managerial habit persistence and time-varying discount rates. However, when discount rates are high (low), elevated dividend-price ratios predict lower (higher) subsequent returns. This result suggests that managerial and consumer habits exert opposing effects on return predictability. Furthermore, two empirical implications derived from the theoretical results are not rejected, underscoring the model's empirical relevance for asset pricing.

CO234 Room Virtual R01 ADVANCEMENT ON STATISTICAL NETWORK MODELING AND ANALYSIS

Chair: Wen Zhou

C1475: **Statistical inference for network regression with random designs**

Presenter: **Yuan Zhang**, The Ohio State University, United States

Network regression incorporates network data into a regression model to predict a nodal response variable. The seminal prior study studied the estimation and inference methods when the design, which refers to both the network model and the predictors, is fixed. The aim is to investigate a different scenario where the design is random. It is observed that a very different source of stochastic variation dominates. Consequently, the statistical inference procedure is very different than that in the prior study. Algorithm, theory, and numerical results are presented for this new problem.

C1476: **Nonparametric estimation for time-varying network models**

Presenter: **Tianxi Li**, University of Minnesota, United States

Co-authors: Adam Rothman, Jeonghwan Lee

The analysis of time-varying networks, in which interactions evolve over time, poses significant statistical challenges. We address the problem of estimating time-varying edge probabilities within the flexible framework of the time-varying graphon model, which posits proper smooth evolutions in both latent node positions and time. We introduce a stagewise smoothing method to estimate a model that is computationally efficient with a uniform convergence guarantee. The method is demonstrated by the modeling evaluation of the US Congress network.

C1482: **Doubly robust alignment for large language models**

Presenter: **Chengchun Shi**, LSE, United Kingdom

Reinforcement learning from human feedback (RLHF) is considered for aligning large language models with human preferences. While RLHF has demonstrated promising results, many algorithms are highly sensitive to misspecifications in the underlying preference model (e.g., the Bradley-Terry model), the reference policy, or the reward function, resulting in undesirable fine-tuning. To address model misspecification, we propose a doubly robust preference optimization algorithm that remains consistent when either the preference model or the reference policy is correctly specified (without requiring both). Our proposal demonstrates superior and more robust performance than state-of-the-art algorithms, both in theory and in practice. The code is available at <https://github.com/DRPO4LLM/DRPO4LLM>.

CO114 Room BCB 206 CFE SESSION: A TRIBUTE TO H. PESARAN IV

Chair: Joachim Schnurbus

C0407: **Fixed-effect estimators under strict exogeneity for averages of unknown heterogeneous effects**

Presenter: **Myoung-jae Lee**, Korea University, Korea, South

"Fixed-effect estimator (FIX)" with a constant treatment effect specification is popular with panel data, but treatment effects are unknowingly heterogeneous in reality, not constant; also, the panel model can be misspecified in other ways. For these problems, the following points are made using nonparametric "causal reduced forms". Firstly, under a strict exogeneity, FIX is consistent with a possibly negative weighted average of heterogeneous treatment effects despite model misspecifications, if a set of restrictive conditions holds; however, negative weights make the estimand non-causal. Secondly, a "modified FIX" is consistent with the same estimand without those restrictive conditions; this is an improvement, but the estimand is still non-causal. Thirdly, the "modified FIX using only the treatment-varying subsample" is consistent with a nonnegatively weighted average effect, which is thus causal. Simulation and empirical studies demonstrate these points.

C0892: **Testing weak dependence and stationarity by universal inference**

Presenter: **Arnab Bhattacharjee**, Heriot-Watt University, United Kingdom

Co-authors: Joseph Paul, Mark Schaffer

A sequential universal inference (UI) test is presented for detecting strong spatial dependence and non-stationarity in dynamic panel data with expanding dimensions N and T . Starting from a general spatial autoregressive (SAR) model with both contemporaneous and lagged spatial effects, an incremental likelihood-ratio is computed between an unrestricted quasi-MLE and the null of zero contemporaneous and lagged spatial coefficients. The cumulative product of these ratios forms a non-negative martingale with unit mean under the null, so Villes' inequality results in an any-time, finite-sample level- α test as the time dimension T grows, without relying on large- T or large- N approximations. The roles of the two dimensions are interchangeable, hence the argument works in both the cross-section and time dimensions. Inverting the same statistic delivers finite-sample confidence sets for the spatial-autoregressive parameters that remain valid in finite samples. Monte-Carlo evidence shows strong size and power performance relative to the Pesaran CD test, a prior study, and conventional asymptotic LR tests, especially when spatial coefficients lie near the unit-root boundary. Robust extensions are outlined for heavy-tailed and heteroskedastic disturbances, and an asymptotic refinement that results in an exact asymptotically level- α test is derived.

C1473: **Testing for structural changes in factor models with short panel data**

Presenter: **Shuyang Sheng**, The Chinese University of Hong Kong, Shenzhen, China

The aim is to investigate methods for testing structural changes in factor models with short panel data, where $N \rightarrow \infty$ and T is fixed. Under the null hypothesis of no structural change in factor loadings, the approach of a prior study is applied to eliminate factor loadings and estimate the parameters by GMM. Building on this idea, a J-test is proposed for detecting structural changes in factor loadings, which has the correct asymptotic

size. The power of the test is also compared with that of the test proposed by a recent study.

CO196 Room BCB 207 ECONOMETRIC FORECASTING
Chair: Robert Kunst
C0660: Optimal combinations of mean square error and directional forecast accuracy for model selection

Presenter: **Robert Kunst**, Institute for Advanced Studies, Austria

Co-authors: Mauro Costantini

Forecasting economic time series often aims at two targets that may be in conflict: Accuracy in the sense of closeness between forecast and realized value, as it is measured by mean square error, and directional forecast error, with interest exclusively focused on ups and downs. In this context, the problem of selecting a forecast model is studied among two assumed candidates. Minimizing a combined loss function that accounts for both targets jointly is considered, using a weighting scheme. In previous work, Monte Carlo analysis under different scenarios has explored the strength of the procedure. Time-homogeneous univariate and vector autoregressions have been considered, but also generating laws that involve thresholds and structural breaks. Some windows have been identified where the weighted combined targets succeed in improving the accuracy criteria. The interaction is investigated between the weight assigned to directional accuracy in the selection criterion and in the evaluation criterion. One may conjecture that a strong weight for directional accuracy (DA) implies a stronger DA performance of the selected model, but this is not necessarily the case. Simulation evidence of realistic designs is presented, where a stronger DA weight implies deteriorating DA performance. In some cases, an optimal weight can be found where DA performance attains a maximum. In empirical applications, such simulations may be helpful in optimizing the DA weight in the criterion function.

C1066: Metal price uncertainty and macroeconomic dynamics

Presenter: **Ines Fortin**, Institute for Advanced Studies, Austria

Co-authors: Jaroslava Hlouskova

A new index measuring metal price uncertainty, following existing methodology, is constructed. The effects of metal price uncertainty shocks are then examined on economic activity as well as on the stock and commodity markets, considering impulse response functions in a vector error correction model setup. In addition, the relative importance of metal price uncertainty is assessed as a leading indicator for real and monetary developments in a forecasting analysis. Concerning commodities, the focus is on industrial metals traded at the London Metal Exchange. Empirical analyses are performed for the Euro area and the United States.

C1459: Macroeconomic forecasting with time series foundation models

Presenter: **Rolf Scheufele**, Swiss National Bank, Switzerland

Time series foundation models (TSFMs) such as TimesFM and Chronos promise major advances in forecasting, yet their evaluation is hampered by limited control over training data and the infeasibility of standard pseudo out-of-sample tests. The purpose is to address this by conducting a real-time forecasting experiment using daily financial and macroeconomic series (exchange rates, interest rates, stock prices) from June 2025 onward. This high-frequency design yields genuine out-of-sample observations, allowing a fair comparison between leading TSFMs and traditional benchmarks. Results show heterogeneous performance: while some TSFMs match or slightly outperform univariate models and the random walk, others fall short. Early evidence is provided on the opportunities and limitations of TSFMs in economics and finance.

CO191 Room BCB 209 MODELLING RISK AND UNCERTAINTY
Chair: Paulo Rodrigues
C1164: High-dimensional panel expectiles regression: A decomposition of the gender wage gap

Presenter: **Pedro Raposo**, Catolica Lisbon school of business and economics, Portugal

Co-authors: Paulo Rodrigues, Pedro Portugal, Matei Demetrescu

The aim is to develop expectile panel data methods for high-dimensional fixed effects estimation in line with a prior study, which allows for a wide range of applications in fields such as Labor economics, Economics of Education, and Inequality. It is shown how the Gelbach decomposition can be validly implemented in the context of panel expectile regressions. Using a unique Portuguese-linked employer-employee dataset, use the estimator to explore the determinants of the gender wage gap over the period 1995-2022. It is found that: (i) the gender wage gap is larger in the upper tail; (ii) the difference is mainly explained both in the left and right tail by the individual unobserved heterogeneity; and (iii) assortative matching is less pronounced in the tails.

C1166: Forecasting the tail behavior of banks risk sentiment

Presenter: **Paulo Rodrigues**, Universidade Nova de Lisboa, Portugal

Co-authors: Joao Nicolau, Adriana Cornea-Madeira

The aim is to provide novel approaches for modelling, forecasting and identifying the drivers of the tail behavior of time series. Using a new triple adaptive lasso approach, the relevance of different macrofinancial proxies are analyzed for forecasting the right tail extreme values of Bank Risk Sentiment Tracker (BRST) series, a synthetic indicator of market sentiment for individual listed banks. Specifically, considering a sample of 21 banks from the Euro area, the UK, and the US, the statistical significance of global and regional market volatility and, skewness measures as well as monetary and macro conditions are evaluated via inflation expectations and the yield curve spread. An extensive and novel forecasting exercise is performed in order to evaluate how well these variables perform in anticipating the risk dynamics traced by the BRST. The focus is on two forecast horizons: One-day ahead, $h = 1$, and five-days ahead $h = 5$.

C1191: If at first you don't succeed: A dynamic evaluation of grade retention

Presenter: **Hugo Reis**, Banco de Portugal, Portugal

Many European countries exhibit high rates of secondary school grade retention. The aim is to investigate the retention's effects on student learning, dropout rates, and educational attainment in Portugal - a country with some of the highest retention rates globally. To analyze these effects, an extended Roy model is developed and estimated that captures the cumulative impact of retention on test scores, dropout, and attainment across multiple grades. The findings reveal substantial heterogeneity in retention's effects on students' educational attainment and academic achievement. A negative consequence of Portugal's high rate of retention (>40%) is that it significantly increases the likelihood of students dropping out of school, by 25 pp for the average retained student. However, for retained students who complete secondary school, 79% experience test score gains in math and 45% experience improvements in Portuguese. To assess model validity, the model is used to simulate retention effects on dropout for 12th-grade students at the margin of being retained and compare the predictions to regression discontinuity (RDD) estimates, which are found to be close. The estimated model is then used to solve for an optimal retention policy that maximizes average lifetime earnings, accounting for retention's countervailing effects on educational attainment and academic achievement.

CO340 Room BCB 210 ADVANCES IN NON-LINEAR TIME SERIES
Chair: Dario Palumbo
C0682: Dynamic combination and calibration for climate predictions

Presenter: **Roberto Casarin**, University Ca' Foscari of Venice, Italy

Co-authors: Dario Palumbo, Francesco Ravazzolo

When multiple forecasts are available from different models or sources, it is possible to combine these to make use of all relevant information on the variable to be predicted and, as a consequence, to produce better forecasts. This is particularly important when working with uncertain

environments, and the selection of relevant information a priori is not an easy task. Climate change and challenges related to global warming are essential issues for science; therefore, modeling their uncertainty and producing reliable forecasts to deal with them are crucial tasks for econometricians. Dynamic combination and calibration are introduced to produce accurate climate predictions. The density calibration literature is extended, and the application of dynamic combinations when calibrating models is proposed. The static model of the prior study is extended to a score-driven dynamic model for calibration and combination of predictive distributions. The time-varying weights are fitted by an observation-driven model with dynamics driven by the score of the assumed conditional likelihood of the data-generating process. The model is very flexible and can handle different shapes, instability, and model uncertainty. It is illustrated analytically and in simulation exercises. Then, the methodology is applied to climate predictions.

C0769: From rotational to scalar invariance: Enhancing identifiability in score-driven factor models

Presenter: **Emilija Dzuverovic**, Ca Foscari University of Venice, Italy

Co-authors: Fulvio Corsi, Giuseppe Bucheri

It is shown that, for a certain class of scaling matrices including the inverse square-root of the conditional Fisher information, score-driven factor models are identifiable up to a multiplicative scalar constant under very mild restrictions. This result has no analog in parameter-driven models, as it exploits the different structure of the score-driven factor dynamics. Consequently, score-driven factor models overcome the issue of rotational invariance that typically affects dynamic factor models, thereby enhancing the economic and financial interpretability of the estimated factors. The restrictions are order-invariant and can be generalized to score-driven factor models with dynamic loadings and nonlinear factor models. The identification strategy is tested extensively, using simulated and real data. The empirical analysis on financial and macroeconomic data reveals a substantial increase in log-likelihood ratios and significantly improved out-of-sample forecast performance when switching from the classical restrictions adopted in the literature to the more flexible specifications.

C0840: Identification, estimation, and inference for models with multiple behavioral equilibria

Presenter: **Davide Raggi**, CaFoscari University of Venice, Italy

Deviations from the rational expectation paradigm pose many challenges in terms of statistical properties of the underlying macroeconomic models and on the theoretical properties of the estimators for the parameters of interest. The focus is on bounded rationality, in which agents use a simple misspecified autoregressive model to build their subjective expectations. This assumption may lead to multiple expectational equilibria. It is also assumed that agents recursively update the parameters of this autoregressive rule. This feature refers to learning, where economic agents update their expectations based on recent historical data. It is found that under constant gain learning, trajectories of the variables of interest are still uniformly ergodic. Furthermore, even under learning, it is possible to derive key properties for the estimators of the structural economic parameters, such as consistency and asymptotic distributions of those estimators.

CO295 Room BCB 211 STATISTICAL MODELS FOR BUSINESS AND ECONOMICS

Chair: Sujay Mukhoti

C0229: Climate-driven inventory optimization using generalized newsvendor model with random supply and demand

Presenter: **Soham Ghosh**, Indian Institute of Technology Indore, India

Co-authors: Sujay Mukhoti, Pritee Sharma

A decision-theoretic framework is developed to optimize humanitarian food supply in drought-prone African countries using an extended newsvendor model that incorporates climate and production risk. Agricultural output is modeled as a function of temperature and latent volatility in precipitation, derived from a Bayesian stochastic volatility structure applied to the autoregressive standardized precipitation index (SPI). MIDAS (mixed data sampling) regressions are employed with beta polynomial weighting to link high-frequency climate risk measures to annual production. The expected production is embedded within a piecewise-linear cost function that penalizes both excess supply and shortfalls. Using particle filtering, the distribution of latent volatility is approximated, and nested expectations that define the cost function are evaluated. The proposed model captures the asymmetric and nonlinear influence of climate risk on food security outcomes and yields optimal inventory decisions under uncertainty. This framework supports adaptive food aid planning by integrating statistical forecasting, volatility modeling, and supply chain optimization. The approach is particularly suited for use by institutions such as the World Food Program, enabling more efficient and responsive resource allocation. The methodology contributes to the intersection of climate econometrics, statistical decision theory, and inventory management under risk.

C0339: Likelihood based estimation of optimal order quantity

Presenter: **Sujay Mukhoti**, Indian Institute of Management Indore, India

Co-authors: Soham Ghosh

Three distinct approaches are proposed and evaluated: Maximum likelihood estimation (MLE), Markov chain Monte Carlo (MCMC), and non-parametric estimation for determining the optimal order quantity in inventory systems under demand uncertainty. As a benchmark, maximum likelihood estimation is employed to fit conventional demand distributions (e.g., log-normal, generalized beta), and the optimal order quantity is computed based on estimated parameters. This method offers analytical tractability and ease of implementation, particularly when demand closely follows known parametric forms. Alternatively, MCMC-based estimators are constructed to account for parameter and model uncertainty by exploring posterior distributions over both space and structure. To provide a more assumption-agnostic alternative, a non-parametric method is also implemented based solely on historical demand data. This approach directly minimizes the empirical cost function over order quantities without imposing any parametric form, thereby enhancing flexibility and robustness in real-world contexts where demand may be skewed. A detailed simulation study compares these three approaches across varying demand scenarios.

C0534: Predictive modeling for coal production locations with kernel density and some machine learning approaches

Presenter: **Karthik Sriram**, Indian Institute of Management India, India

For a major coal mining company in India, the problem of predicting which of its mining locations will produce transportable coal on any given day is addressed. The performance of different predictive machine learning approaches is adapted and evaluated, including random forests, support vector machines, logistic regression, linear discriminant analysis, multilayer perceptrons, and the long short-term memory model. In addition, motivated by the context and structure of the coal production data, a novel approach is devised based on kernel density estimation (KDE) with adaptive bandwidths, by accounting for the discrete spatial locations of the mines as well as their varying frequencies of production. The different modeling approaches are evaluated using standard predictive accuracy measures for classification problems, based on recent production data from 136 different mines (owned by the coal mining company) over one year. The results suggest that the KDE-based approach outperforms standard machine learning approaches, highlighting the importance of devising novel modeling solutions motivated by the structure of contextual data.

C1058: High-dimensional regularized additive matrix autoregressive model

Presenter: **Nilanjana Chakraborty**, Indian Institute Of Management Udaipur, India

Co-authors: Debika Ghosh, Samrat Roy

High-dimensional time series have diverse applications in econometrics and finance. Recent models for capturing temporal dependence have employed a bilinear representation for matrix time series, or the Tucker-decomposition-based representation in the case of tensor time series. A bilinear or Tucker-decomposition-based temporal effect is difficult to interpret on many occasions, along with its computational complexity due to the non-convex nature of the underlying optimization problem. Moreover, the existing matrix case models have not sufficiently explored the possibilities of imposing any lower-dimensional pattern on the transition matrices. A regularized additive matrix autoregressive model is proposed

with additive interaction of row-wise and column-wise temporal dependence, which offers more interpretability, less computational burden due to its convex nature, and estimation of the underlying low rank plus sparse pattern of its transition matrices. The issue of identifiability of the various components in the model is addressed, and subsequently, a scalable alternating block minimization algorithm is developed for estimating the parameters. A finite sample error bound is provided under high-dimensional scaling for the model parameters. Finally, the efficacy of the proposed model is demonstrated on synthetic and real data.

CO127 Room BCB 212 ECONOMIC UNCERTAINTY AND ITS EFFECTS
Chair: Svetlana Makarova
C0985: A new measure of geopolitical risk for a large cross-section of countries

Presenter: **Karol Szafranek**, SGH Warsaw School of Economics, Poland

Co-authors: Joscha Beckmann, Michal Rubaszek, Michael Murach

The aim is to propose new geopolitical risk (GPR) indices for over 180 countries and the global economy. A unique database on activity metrics, sentiment, and tonality on geopolitical events implied from newspapers and social media is explored via natural language processing. Narrow and broad GPR indices are provided since 1998, and the measures are compared to the popular benchmark presented in a prior seminal work. It is shown that the narrow GPR index for the global economy is strongly correlated with the CI measure, which is not the case for the broad GPR proxy. Next, evidence is provided that in specifications using the broad measure, the response of key macroeconomic variables to GPR changes is stronger, while the portion of variance explained by the model increases. It is also found that negative media sentiment acts as a propagation mechanism for geopolitical risk, in particular in high-income and OECD economies.

C1115: Rethinking GPR: The sources of geopolitical risk

Presenter: **Javier J Perez**, Bank of Spain, Spain

Co-authors: Marina Diakonova, Irma Alonso

Geopolitical risks are increasingly cited by policymakers and analysts as key factors influencing economic activity in both the short and medium term. A widely used method to measure these risks involves counting news articles on adverse geopolitical events, notably following a prior study's geopolitical risk (GPR) indexes. The aim is to propose a new approach to enhance GPR measurement by decomposing the index based on the origin of the risk. The concept of bilateral GPRs is introduced, which links geopolitical risks to specific countries or entities. Aggregating these bilateral indexes offers a clearer interpretation of overall GPR. Findings show that these new indices provide distinct insights beyond the benchmark GPR, offering a more accurate picture of current global tensions. Moreover, it is demonstrated that the geographical source of a GPR shock significantly affects its macroeconomic impact both in magnitude and direction on a given economy, as shown through standard VAR models.

C1339: Testing for the footprints of stabilization economic policy in forecast errors

Presenter: **Svetlana Makarova**, University College London, United Kingdom

Co-authors: Wojciech Charemza, Christian Francq, Radu Lupu, Jean-Michel Zakoian

The purpose is to introduce a novel statistical test, the policy effects Lagrange multiplier (PELM) test, to detect stabilization policy effects in the distribution of forecast errors from dynamic financial models. Traditional analyses of policy impact typically rely on explicit policy information or direct intervention data, which are often unavailable or incomplete. In contrast, the proposed PELM test infers policy footprints from the distribution of forecast errors alone. Empirically applied to sovereign bond yield data from 33 countries before the Russian financial crisis of 2014, the test identifies countries showing stabilization policy footprints. Subsequent analysis shows that significant budgetary improvements were observed for years following the crisis in the group of countries where the test statistically confirmed stabilization policies. This confirms the rationale of test foundations and also indicates its predictive properties. Robustness checks further validate these findings across various model specifications and sensitivity scenarios. The proposed PELM test offers policymakers and researchers a powerful tool for evaluating stabilization policies, facilitating better forecasting and assessing policy efficiency in diverse economic contexts without necessitating detailed policy intervention data.

CO338 Room BCB 213 THE ECONOMETRICS OF EMPIRICAL ASSET PRICING
Chair: Andre B M Souza
C1251: Multivariate factors: Accounting for the joint dependence among characteristics

Presenter: **Anastasija Teterova**, Erasmus University Rotterdam, Netherlands

Co-authors: Rasmus Lonn, Gustavo Freire

The aim is to propose a new approach for constructing characteristics-based factor portfolios. Instead of forming independent quantiles or univariate rank sorts, weights are allocated proportional to the conditional distribution of each characteristic given all other characteristics. This is modelled as a conditional distribution in a simple and data-driven way using copulas. The method is applied to the five Fama-French factor portfolios. Relative to the original construction, the multivariate factors increase the maximum attainable Sharpe ratio from 1.04 to 1.80 and decrease by half the number of anomalies with significant alphas in the cross-section of stock returns.

C1492: Transaction costs and the stochastic discount factor

Presenter: **Daniele Bianchi**, Queen Mary University of London, United Kingdom

Transaction costs are shown to fundamentally reshape the stochastic discount factor (SDF) by determining which characteristics matter in equilibrium. Using deep neural networks, transaction costs are incorporated directly into robust SDF estimation rather than treated as post-optimization adjustments or arbitrary investment constraints. Transaction cost-aware SDFs yield substantially higher net Sharpe ratios and superior cross-sectional pricing through endogenous portfolio reallocation, including increased diversification, reduced turnover, and lower exposure to costly high-turnover characteristics. These effects persist across sample configurations, market regimes, neural network specifications, and alternative definitions of transaction costs, demonstrating that trading frictions are structural determinants of equilibrium asset prices.

C1497: How to bet on winners (and losers)

Presenter: **Andre B M Souza**, ESADE Business School, Spain

Co-authors: Christian Brownlees

The construction of long-short portfolios is cast as a statistical decision problem in which the investor seeks to buy the top-performing stocks (the "winners") and sell the worst-performing ones (the "losers") on the basis of stock characteristics. We derive the optimal portfolio selection rule implied by a loss function that accounts for different types of misclassification errors in portfolio construction. This approach leads to a return classification problem and the optimal rule buys or sells stocks based on their probabilities of being winners or losers, conditional on the stock characteristics. When returns are generated by an additive regression model and misclassification costs satisfy a symmetry condition, the optimal rule simplifies to the conventional sorting procedure based on expected returns. An empirical application using U.S. stock data shows that portfolios constructed using the optimal rule achieve higher Sharpe ratios compared to those built using conventional methods. Our results demonstrate that predictive signals in the cross-section of stock returns go beyond expected returns, and that properly optimized portfolio selection rules based on these signals can generate substantial economic value for investors.

CO106 Room BCB M201 EXTREMES AND DEPENDENCE MODELLING
Chair: Boris Beranger
C0270: A class of skew-multivariate distributions for spatial data

Presenter: **Pavel Krupskiy**, Melbourne University, Australia

A class of copula models is introduced for spatial data based on multivariate Pareto-mixture distributions. The tail properties of these models are explored, demonstrating their ability to capture both tail dependence and asymptotic independence, as well as the tail asymmetry frequently observed in real-world data. The proposed models also offer flexibility in accounting for permutation asymmetry and can effectively represent both the bulk and extreme tails of the distribution. Special cases of these models are considered with computationally tractable likelihoods, and an extensive simulation study is presented to assess the finite-sample performance of the maximum likelihood estimators. Finally, the models are applied to analyze a temperature dataset, showcasing their practical utility.

C0777: On factor copula-based mixed regression models

Presenter: **Bouchra Nasri**, University of Montreal, Canada

Co-authors: Bruno Remillard, Pavel Krupskiy

A copula-based method for mixed regression models is proposed, where the conditional distribution of the response variable, given covariates, is modeled by a parametric family of continuous or discrete distributions, and a latent variable models the dependence between observations in each cluster. The estimation of copula and margin parameters is demonstrated, outlining the procedure for determining the asymptotic behavior of the estimation errors. Numerical experiments are performed to assess the precision of the estimators for finite samples. An example of its application is given using COVID-19 vaccination hesitancy data from several countries. All developed methodologies are implemented in CopulaGAMM, available in CRAN.

C0949: Statistical methods for multivariate time series with arbitrary distributions

Presenter: **Bruno Remillard**, HEC Montreal, Canada

Co-authors: Bouchra Nasri, Kilani Ghoudi

The aim is to present recent statistical methods that can be used for multivariate time series with arbitrary distributions, i.e., the associated distributions can be continuous, discrete, or a mixture of continuous and discrete, such as zero-inflated distributions. Using a new transformation, it is possible to define tests of independence between time series, as well as tests of goodness-of-fit.

CO184 Room BCB M202 BAYESIAN METHODS IN FINANCE

Chair: Maria Fernanda Pintado

C0203: Intraday crude oil volatility: Assessing the impact of economic announcements and mixed-frequency data

Presenter: **Audrone Virbickaite**, CUNEF Universidad, S.L., Spain

Co-authors: Igor Ferreira Batista Martins, Hoang Nguyen, Hedibert Lopes

The real-time impact of economic announcements and policy decisions is investigated on intraday crude oil volatility. The model employs a high-frequency mixed data sampling (MIDAS) approach to capture the interaction between market volatility and economic uncertainty, paired with spike-and-slab priors for the announcement coefficients to perform efficient variable selection in a high-dimensional setting. Findings demonstrate that these announcements significantly influence short-term oil price volatility; Meanwhile, the evolution of the level of the volatility depends on some economic indicators, sampled at lower frequencies. The proposed model improves the accuracy of short-term volatility forecasts, offering valuable insights for market participants and policymakers. The empirical results highlight the importance of timely economic information in forecasting oil market dynamics.

C0880: Flexible combination strategies for multivariate volatility forecasting

Presenter: **Christoph Frey**, Lancaster University, Germany

Co-authors: David Happersberger

Accurately forecasting multivariate volatility is crucial for risk management and portfolio allocation, yet model uncertainty and the curse of dimensionality often hinder the identification of a single optimal forecasting model. Additionally, it is well-documented in the literature that model combinations, across various domains, can significantly improve the accuracy of out-of-sample forecasts. The aim is to propose novel combination strategies that extend existing approaches by integrating advanced weighting schemes and dynamic model selection to improve the predictive accuracy of multivariate volatility forecasts. The methodology leverages multivariate conditional volatility models, realized covariance matrices, and dynamic-factor models, combining them with model weights that adjust based on recent individual forecast performance or portfolio utility. Comprehensive empirical analysis demonstrates that the proposed combination methods consistently outperform individual models and traditional equal-weighted combinations in terms of statistical accuracy and economic relevance, while also reducing parameter dependence.

C1344: Bayesian Markov-switching partial reduced-rank regression

Presenter: **Maria Fernanda Pintado**, CUNEF Universidad, Spain

Co-authors: Matteo Iacopini, Luca Rossini, Alexander Shestopaloff

Reduced-rank (RR) regression is a powerful dimensionality reduction technique, but traditional RR models typically overlook any potential group structure among the responses by assuming a low-rank structure on the coefficient matrix. When the observations in the regression model are indexed by time, the relationship between covariates and responses could change over periods. A time-varying grouping structure in the response variables in RR regression is currently understudied. To address this limitation, a Markov-switching Bayesian partial RR (MSPRR) regression is proposed. First, the response vector is partitioned into two groups to reflect different degrees of complexity of the relationship. A "simple" group assumes a low-rank linear regression, and a "flexible" group remains agnostic and exploits nonparametric regression via a Gaussian process. Second, different from traditional approaches that assume known group structure and rank, they are treated as unknown parameters to be estimated. Third, time variation and persistence are accounted for by introducing a Markov-switching process, which examines the changes in the grouping structure and model parameters over time. Fully Bayesian inference is performed via a partially collapsed Gibbs sampler, which allows uncertainty quantification. Applications to both synthetic and macroeconomic data demonstrate the capability of the proposed method to uncover latent states and hidden structures within the data.

CO041 Room BCB 308 OVER-PARAMETRIZATION AND OVERFITTING IN MACHINE LEARNING

Chair: Debarghya Ghoshdastidar

C0414: Implicit regularization of deep residual networks towards neural ODEs

Presenter: **Pierre Marion**, Inria, France

Residual neural networks are state-of-the-art deep learning models. Their continuous-depth analog, neural ordinary differential equations (ODEs), are also widely used. Despite their success, the link between the discrete and continuous models still lacks a solid mathematical foundation. A convergence result of deep residual networks towards neural ODEs is discussed, for nonlinear networks trained with gradient flow: If the network is initialized as a discretization of a neural ODE, then such a discretization holds throughout training. The result is valid for a finite training time, and also as the training time tends to infinity, provided that the network satisfies a Polyak-Lojasiewicz condition. Importantly, this condition holds for a family of residual networks where the residuals are two-layer perceptrons with an overparameterization in width that is only linear, and implies the convergence of gradient flow to a global minimum. If time allows, consequences will be also discussed in terms of statistical guarantees, namely generalization bounds.

C1049: Kernel regime of deep neural networks: Insights and limitations

Presenter: **Mariia Seleznova**, Ludwig Maximilian University of Munich, Germany

Training dynamics of non-linear deep neural networks (DNNs) are famously challenging to analyze, so current theory heavily relies on simplifica-

tions. In particular, DNN dynamics simplify dramatically in the infinite-width limit, entering the so-called "kernel regime" under certain conditions. In this regime, the dynamics are linearized around the initialization and governed by a deterministic and constant neural tangent kernel (NTK). This allows for a theoretical treatment of optimization and generalization using the NTK - an approach adopted in many recent studies. Given that modern DNNs are typically overparameterized, the infinite-width limit appears to offer a promising framework for these models. However, several limitations of this approach have been identified, and it is examined whether the kernel regime indeed provides a good approximation for the behavior of deep fully connected networks. The results reveal that the depth-to-width ratio and the initialization distribution play a critical role. In particular, very deep networks are generally not in the kernel regime at the beginning of training. A new approach is also proposed to study the dynamics using the "kernel regime at the end of training", which enables the prediction of the neural collapse (NC) phenomenon.

C0528: Self-supervised learning in the kernel regime

Presenter: **Maximilian Fleissner**, Technical University of Munich, Germany

In recent years, self-supervised learning (SSL) has emerged as a powerful paradigm, building the foundation of several modern machine learning models. At its core, SSL relies on the idea of using data augmentations to encode a notion of similarity in otherwise unlabeled samples. However, despite its rapidly growing popularity, the statistical understanding of what SSL learns is still limited. Arguably, one of the most promising avenues towards understanding the fundamental principles of SSL is by connecting it to kernel methods. In supervised learning, this correspondence is justified by virtue of the neural tangent kernel (NTK), which asserts that overparameterized neural networks follow training dynamics that resemble those of kernel machines. An extension of the NTK to a commonly used SSL algorithm is presented, namely, Barlow Twins. The NTK connection allows characterizing the patterns learned in SSL, as well as quantifying their generalization properties.

CO206 Room BCB 309 SAFE AI FOR DECISION-MAKING IN ECONOMICS AND FINANCE (VIRTUAL)

Chair: Emanuela Raffinetti

C0182: Focused SAFE-driven AI approach for volatility prediction in EU ETS

Presenter: **Emanuela Raffinetti**, University of Pavia, Italy

Co-authors: Maria Elena De Giuli

The EU emissions trading system (EU ETS) involves regulatory changes and sectoral transformations whose non-linear dynamics require the employment of powerful forecasting techniques. Recent literature showed that artificial intelligence (AI) systems appear more effective than classical econometric approaches in detecting the complexities of the EU ETS. Despite their advantages, the black-box nature of AI systems may give rise to direct outputs without a clear connection with the inputs that generate them. This is the reason why highly complex machine and deep learning models have to be supported by actions addressed to avoid, or at least restrict, the dangerous effects derived from their incorrect use. To this purpose, the European Commission has promoted focused safety principles for a trustworthy AI, specified in terms of sustainability, accuracy, fairness, and explainability. In line with these premises, the contribution is twofold: 1) on the theoretical side, the metrics associated with the principles of accuracy and explainability is formalized in order to investigate the role of historical data in shaping the realized volatility of the EUA prices and the implied volatility inferred from the EUA-related financial instruments; 2) on the empirical side, the proposed methodology is implemented on EUA price data.

C0734: Matrix completion for spatiotemporal air quality datasets

Presenter: **Rodolfo Metulini**, University of Bergamo, Italy

Matrix completion (MC) is a recent and flexible statistical learning method for imputing missing values and performing counterfactual analysis in structured datasets. A key advantage of MC methods is that they avoid relying on stringent model assumptions. A widely used approach leverages nuclear norm regularization, which promotes low-rank approximations of the data matrix. Recent advances on nuclear norm MC incorporate unit- and time-specific fixed effects, which is crucial to avoid biased imputations in panel data where strong heterogeneity across units and time is expected. The performance of existing MC methods is evaluated on panel data from the ARPA Lombardia air quality monitoring network. Given the spatial correlation inherent in air pollution data, the aim is also to incorporate spatial constraints into the matrix completion framework and to evaluate model accuracy and interpretability. The development and testing of MC methods contribute to the broader goal of conducting accurate counterfactual analyses in policy evaluation, particularly for assessing the impact of mobility-related interventions on air pollution levels. By enabling reliable estimations of air quality in the absence of specific policies, such methods can support informed decision-making in environmental and urban planning.

C1517: Non-parametric causal discovery for EU allowances returns through the information imbalance

Presenter: **Cristiano Salvagnin**, University of Insubria, Italy

Co-authors: Aldo Glielmo, Vittorio Del Totto, Antonietta Mira, Maria Elena De Giuli

A new non-parametric method called Differentiable Information Imbalance is applied to identify variables that are causally linked, possibly through complex and non-linear relationships, to the financial returns of European Union Allowances in the European Union Emissions Trading System. Using data from January 2013 to April 2024, we compare this approach with the traditional multivariate Granger causality method based on vector autoregressive models. Both methods identify key drivers such as coal futures prices and the Spanish stock index IBEX 35, but they also reveal notable differences. Through experiments with synthetic data, we show that these differences likely arise from the linear assumptions underlying standard Granger causality.

CO315 Room BCB 310 INTERACTING URNS AND INNOVATION PROCESSES

Chair: Andrea Ghiglietti

C0324: A system of urn models for incorporating informational borrowing in the design and inference of clinical trials

Presenter: **Rosamarie Frieri**, University of Bologna, Italy

Co-authors: Andrea Ghiglietti, Giacomo Aletti, Irene Crimaldi, Alessandro Baldi Antognini

A new design methodology is introduced for stratified comparative experiments based on a system of interacting urns. The key idea is to model the interaction between urns for borrowing information across strata and to use it in the design phase in order to i) enhance the information exchange at the beginning of the study, when only a few subjects have been enrolled and the stratum-specific information on treatments efficacy could be scarce, ii) let the information sharing adaptively evolves via an update mechanism based on the observed outcomes, for skewing at each step the allocations towards the stratum-specific most promising treatment and iii) make the contribution of the strata with different treatment efficacy vanish as the stratum information grows. In particular, the interacting urns design is introduced, namely a new covariate-adjusted response-adaptive procedure that randomizes the treatment allocations according to the evolution of the urn system. The theoretical properties of this proposal are described, and the corresponding asymptotic inference is provided. Moreover, by a functional central limit theorem, the asymptotic joint distribution of the Wald-type sequential test statistics is obtained, which allows for sequential monitoring of the suggested design in clinical practice.

C0804: Generalized measure-valued Polya sequences

Presenter: **Hristo Sariev**, Sofia University, Bulgaria

Measure-valued Polya urn sequences (MVPS) are a generalization of the observation processes generated by k-color Polya urn models, where the space of colors is unbounded and urn composition is modeled by finite measures. An extension of MVPSs is investigated via a randomization of the law of the reinforcement, called generalized measure-valued Polya urn sequences (GMVPS). Using that, the urn composition process is a measure-valued Markov chain, showing that GMVPSs can be represented as mixtures of MVPSSs. The focus is then on the class of generalized randomly reinforced Polya sequences (GRRPS), which are GMVPSs whose reinforcement is a weighted Dirac measure. It follows from the

form of their predictive distributions that the dynamics of GRRPSs are driven by the interaction between weights and observations. Under the additional assumption that the weights are marginally exchangeable, it is proven that the joint process tracking weights and observations is partially conditionally identically distributed, from which its asymptotic properties are derived.

C1072: Modeling innovation ecosystem dynamics through interacting reinforced Bernoulli processes.

Presenter: **Federico Nutarelli**, IMT Lucca, Italy

Co-authors: Irene Crimaldi, Andrea Ghiglietti, Giacomo Aletri

Understanding how capabilities grow into core strengths, and how those strengths can harden into rigidities, is central to innovation strategy. Yet formal modeling is difficult because innovations in one area often reshape others, making specialization endogenous. This challenge is acute in ecosystems where firm performance depends on managing interdependencies and complementarities. The aim is to introduce a formal model based on interacting reinforced Bernoulli processes that tracks how patent wins propagate across technology categories and how those categories co-evolve. The model reproduces key empirical regularities: sublinear growth in cumulative success, convergence of success shares across fields, and declining cross-category correlations over time. Using GLOBAL PATSTAT data (1980-2018), the model is validated, the structural interaction matrix is estimated, and a statistical test is developed for the strength of cross-category effects under a mean-field approximation. By endogenizing specialization, the framework offers a practical tool for policymakers and managers who must steer complex, co-evolving innovation systems.

C1521: Lost in the supermarket: Individual and collective dynamics through the lens of innovation processes

Presenter: **Margherita Lalli**, Scuola Normale Superiore di Pisa, Italy

Co-authors: Francesca Tria, Luca Pappalardo

The focus is on an urn-based model aimed at characterizing the interplay between individual and collective exploration in a population of interacting agents. The model is inspired by the feedback dynamics of recommender systems, where users' choices and algorithmic retraining form a continuous loop that can generate unintended emergent effects. A central question, widely debated in computer science, is whether and under which conditions increasing individual-level diversity in choices may paradoxically reduce diversity at the collective scale. Conclusive results are particularly hindered by the opacity of real recommendation algorithms, which motivates the need for generative models capable of disentangling the underlying mechanisms. Our framework is empirically guided by a rich dataset of supermarket transactions featuring statistical regularities typical of innovation processes such as Zipf's, Heaps', and Taylor's laws, alongside systematic discrepancies between individual and aggregate purchases. To capture this, we build upon the Urn Model with Triggering, an extension of Polya's urn capable of reproducing such phenomenologies at a single scale. Unlike previous models of interacting urns, our formulation assumes neither a predefined network structure nor a nested user-product organization, offering a complementary perspective on the mechanisms driving the tension between individual exploration and collective regularities in consumption.

CO276 Room BCB 311 UNIFORM INFERENCE ON HIGH-DIMENSIONAL MODELS

Chair: Takahiro Nishiyama

C1324: Estimation of latent group structures in time-varying panel data models

Presenter: **Paul Haimerl**, Aarhus University, Denmark

Co-authors: Ines Wilms, Stephan Smeekes

The purpose is to introduce a panel data model where coefficients vary both over time and over cross-section. Slope coefficients change smoothly over time and follow a latent group structure, being homogeneous within but heterogeneous across groups. The group structure is identified using a pairwise adaptive group fused-Lasso penalty. The trajectories of time-varying coefficients are estimated via polynomial spline functions. The asymptotic distributions of the penalized and post-selection estimators are derived, and their oracle efficiency is shown. A simulation study demonstrates excellent finite sample properties. An application to the emission intensity of GDP highlights the relevance of addressing cross-sectional heterogeneity and time variance in empirical settings.

C0607: Inference and automatic instrument selection in a semiparametric single-index conditional factor model

Presenter: **Qi Zhang**, Southwestern University of Finance and Economics, China

Co-authors: Qihua Xu, Chen Huang

A semi-parametric conditional latent factor model with a single-index structure is investigated for estimation and variable selection. The sieve method and singular value decomposition are employed to estimate latent factors, while adaptive LASSO is utilized to identify relevant characteristics. Both large-N and large-NT asymptotics are established for all estimators and the pricing-error test. Variable selection exhibits consistency and Oracle properties. A weighted bootstrap procedure is developed to test the null hypothesis that pricing errors are zero. Empirical analysis of U.S. equity data illustrates out-of-sample predictive performance and the set of selected characteristics.

C1193: Testing block-diagonal covariance structure under a low dimensional factor model in high-dimensional settings

Presenter: **Takahiro Nishiyama**, Senshu University, Japan

Co-authors: Masashi Hyodo, Shoichi Narita

The aim is to propose a new test for block-diagonal covariance structure in a high-dimensional framework, while accommodating a low-dimensional latent factor model. The test, built under low-dimensional factor models, distinguishes from previous normal approximation-based tests, which are valid under a weak spike structure. A modified RV coefficient is proposed for high-dimensional data, and it is shown that its null-limiting distributions follow a weighted mixture of chi-square distributions under a high-dimensional asymptotic regime integrated with weak technical conditions. By applying this asymptotic result and estimation theory of the number of factors in a low-dimensional factor model, a new approximation test is proposed for a block-diagonal covariance structure. The finite sample and dimensional performance of this test are also examined using Monte Carlo simulations.

CO322 Room BCB 312 STATISTICAL METHODS FOR ENVIRONMENTAL SUSTAINABILITY

Chair: Alessandro Bitetto

C0785: Environmental performance and market risk: A data-driven approach

Presenter: **Arianna Agosto**, University of Pavia, Italy

Co-authors: Paola Cerchiello, Alessandra Tanda

The aim is to investigate how environmental performance and climate risk influence market returns and systemic risk. The data-driven environmental performance index (DDEPI) is introduced, a country-level metric constructed using robust principal component analysis (RobPCA). Unlike traditional indices with expert-defined weights, DDEPI is built on empirical data, focusing on variables that account for the largest share of variance, dynamically reflecting variable interrelations across countries and over time. The predictive power of DDEPI is assessed on asset return dynamics, conditional on firm-level ESG ratings and financial indicators. Furthermore, using systemic risk measures such as conditional value-at-risk (CoVaR), it is explored how environmental and climate risks exacerbate firms' exposure to systemic financial shocks. To incorporate physical climate risk, the models are enriched with a novel metric derived from the application of multivariate count time series models that captures the persistence and interdependence of climate-related shocks across sectors and regions. Findings offer new insights into the financial implications of sustainability and environmental issues in market risk assessment.

C0756: Quantitative robustness assessment of composite indicators: Methodology and application

Presenter: **Viet Duong Nguyen**, University Carlo Cattaneo - LIUC, Italy

Benchmarking exercises based on composite indicators are widely used in decision-making across diverse domains, ranging from policy formulation to performance evaluation. However, producing a reliable benchmark poses significant challenges in multi-criteria decision analysis, where outcomes may vary substantially due to minor changes in methods or input features. This uncertainty impacts both developers who strive to generate trustworthy comparisons and decision-makers who rely on these results for strategic planning. As a remedy, a novel approach is introduced for quantifying the robustness of composite indices, accounting for multiple sources of uncertainty such as normalization techniques, weighting schemes, and aggregation functions. The method offers a comprehensive evaluation of index performance based on a model configuration grid, designed to identify the interacting factors that significantly compromise pairwise ranking stability. Its utility is demonstrated through a measuring application using the environmental sustainability dashboard (UNDP), where the proposed method provides valuable insights into benchmarking consistency across computational schemes and highlights critical points for robustness improvement.

C0597: Unpacking determinants of water conservation behavior and the spatial variations of their effects

Presenter: **Rashad Mammadli**, LIUC University C.Cattaneo, Italy

Co-authors: Chiara Gigliarano

Water scarcity issues worldwide demand a thorough understanding of the factors influencing water consumption and conservation behavior to effectively address demand-side challenges. The aim is to investigate the determinants of water conservation behavior through a comprehensive empirical analysis, emphasizing spatial heterogeneity at the regional scale in Italy. Using ordinal logistic regression with sequential analysis, it is found that various socio-demographic, behavioral, and social factors, including trust in public institutions, significantly impact water conservation. Furthermore, applying the mixed geographically weighted spatial autoregressive ordinal regression (MGW-SAR-Ordinal) model reveals statistically significant spatial variation in the relationships between water saving and four key factors: Gender, household size, energy saving, and trust in public institutions. Three of these factors show consistently positive associations with water conservation, while household size exhibits both positive and negative effects depending on the region. These findings highlight the importance of region-specific interventions to promote effective water-saving behaviors, offering valuable insights for policymakers aiming to foster sustainable water conservation practices.

CO198 Room BCB 403 AI FOR COMPLEX NETWORKS: APPLICATIONS IN HEALTH AND FINANCIAL SYSTEMS Chair: Giorgia Riveccio

C0908: Comorbidity patterns and temporal associations of multiple long-term conditions in adults with intellectual disability

Presenter: **Georgina Cosma**, Loughborough University, United Kingdom

Multiple long-term conditions (MLTCs), defined as the presence of two or more chronic health conditions, pose significant challenges for healthcare systems, particularly for vulnerable populations such as individuals with intellectual disability. The distribution, temporal associations, and age/sex-specific patterns of MLTCs are statistically examined in adults with intellectual disability using longitudinal healthcare data. The 18,144 adults with intellectual disability (10,168 males and 7,976 females) are analysed and identified in the Clinical Practice Research Datalink, linked to Hospital Episode Statistics data (2000-2021). Temporal statistical analysis established directional associations among 40 long-term conditions, stratified by sex and age groups (under 45, 45-64, 65+). The high prevalence of enduring mental illness across all age groups distinguished this population from the general population. In males, mental illness co-occurred with gastrointestinal conditions, while in females, mental illness presented alongside chronic pain and endocrine conditions. Findings highlight complex sex-specific MLTC patterns across age groups, revealing temporal associations that provide insights into disease progression, informing targeted prevention strategies and interventions to prevent premature mortality in this vulnerable population.

C0902: Joint modeling of temperature mean and volatility for weather derivative evaluation: A neural network approach

Presenter: **Zelda Marino**, University of Naples Parthenope, Italy

Co-authors: stefania Corsaro, Salvatore Scognamiglio, Vincenzo Di Sauro

The impact of weather on business is significant. Many sectors, such as tourism and agriculture, face weather-related risks, highlighting the need for hedging. Weather derivatives offer an effective tool to reduce the impact of weather variability on operations and financial outcomes. These instruments are based on weather indices over time whose payoffs can depend on variables like temperature, rainfall, or snowfall. Most traded derivatives use temperature indices such as heating and cooling degree days (HDD, CDD), which reflect deviations from a base temperature. Thus, accurate temperature modeling is essential for pricing and risk management. The aim is to propose a neural network approach for calibrating temperature models in weather derivatives. Building on prior work using fully connected networks to estimate daily means, a subnet for variance is introduced. This joint model offers better insight into temperature uncertainty. Assuming a normal distribution, the mean and variance are jointly calibrated. For explainability, the output is designed to mirror known models. Using NASA's MERRA-2 data, it is shown how the model effectively captures seasonal fluctuations in temperature variability. Numerical comparison with existing approaches in the literature shows the effectiveness of the proposed method in terms of both point prediction accuracy and the coverage probability of interval predictions.

C0850: Multilayer perceptron for the economic analysis of diagnostic inappropriateness in cancer networks

Presenter: **Anna Pia Di Iorio**, University of Naples Parthenope, Italy

Co-authors: Giorgia Riveccio, Sandro Pignata, Francesco Schiavone

Innovation is pivotal for optimizing patient care pathways in healthcare, yet its integration into clinical practice demands models that effectively incorporate new technologies. Regional cancer networks serve as strategic structures to coordinate care and promote diagnostic appropriateness, key to cost control and sustainability of health systems. In this context, machine learning-based predictive models are increasingly utilized. The aim is to present a model designed to estimate costs related to diagnostic inappropriateness in oncology, defined as tests not aligned with clinical guidelines or unnecessarily repeated. A multi-layer perceptron (MLP) artificial neural network was used to capture nonlinear relationships among clinical, organizational, and sociodemographic variables. The model was trained on a structured dataset including patients' sociodemographic profiles and oncological care journeys. The target variable was the proportion of inappropriate diagnostic costs per patient. To ensure explainability, SHAP (SHapley Additive exPlanations) analysis was applied to identify key predictors. The most significant factor was time to first multidisciplinary team evaluation, followed by diagnostic modality, distance from the treatment center, and patient age. The results support the development of targeted interventions to reduce inefficiencies and improve diagnostic appropriateness in oncology, contributing to more sustainable and patient-centered care.

CO328 Room BCB 407 RECENT ADVANCES OF REINFORCEMENT LEARNING AND DYNAMIC DECISION-MAKING Chair: Yifan Cui

C1121: On multiple robustness of proximal dynamic treatment regimes

Presenter: **Yuanshan Gao**, Center for Data Science, Zhejiang University, China

Co-authors: Yang Bai, Yifan Cui

Dynamic treatment regimes (DTRs) are sequential decision rules that adapt treatment according to individual time-varying characteristics and outcomes to achieve optimal effects, with applications in precision medicine, personalized recommendations, and dynamic marketing. Estimating optimal DTRs via sequential randomized trials might face costly and ethical hurdles, often necessitating the use of historical observational data. The purpose is to utilize a proximal causal inference framework for the identification of optimal DTRs when the unconfoundedness assumption fails. Contributions are three-fold: (i) proposing a new non-parametric identification method for optimal DTRs with a reduced risk of error amplification; (ii) establishing the semi-parametric theory, including efficient bounds, for the value function of a given sequence of treatment rules; and (iii)

proposing a multiply robust method for identifying and estimating optimal DTRs and proving the corresponding theoretical properties such as consistency and convergence rate. Numerical experiments validate the efficiency and multiple robustness of the proposed methods.

C1425: **Consistent order determination of Markov decision process**

Presenter: **Chuyun Ye**, Beijing Normal University, China

Co-authors: Ruoqing Zhu, Lixing Zhu

Reinforcement learning (RL) leverages the Markov decision process (MDP), which fundamentally relies on the Markov property. However, numerous real-world systems exhibit extended temporal dependencies, demanding higher-order Markov models beyond the typical first-order assumption. The aim is to tackle the challenge of consistently estimating the order of such Markov processes, a problem where traditional sequential testing methods are hindered by limitations in sensitivity and consistency. The purpose is to introduce a novel, two-stage estimation procedure: first, a function is defined that precisely captures the k -order Markov assumption, guaranteeing sensitivity to all violations; second, a signal statistic is constructed that consistently identifies the true order by exploiting a distinct pattern of minimizers. This approach yields a consistent estimator and facilitates efficient implementation. Furthermore, the characteristic curve pattern of the signal statistic aids in visual inspection, that could simplify the order determination process in practical applications. The effectiveness of the method is validated through simulations and a real-world dataset, representing a significant stride in accurately modeling and applying RL to systems with complex temporal dependencies.

C0837: **Reinforcement learning for estimation and inference in heterogeneous environments**

Presenter: **Atanas Christev**, Heriot-Watt University, United Kingdom

Sequential decision-making is central to many social and economic phenomena. Reinforcement learning (RL) is an optimization framework that endows agents with sequential knowledge to interact and learn from an unknown environment, i.e., learn optimal policies of their own behavior in multiple steps through exploration and exploitation. Classical methods struggle to deal with the inherent heterogeneity of structural economic models. Based on inverse reinforcement learning (IRL), a novel estimation framework is proposed and developed that simultaneously clusters behavioral trajectories and infers distinct type-specific utility functions for latent groups, thereby allowing for unobserved heterogeneity. The method supports scalable estimation without explicit transition models and flexible, non-parametric rewards. Theoretically, it is shown that it exhibits oracle properties for group recovery under sufficient reward separation and asymptotic normality, enabling the construction of confidence intervals for policy parameters. The method is applied to analyze the structural representation of a frictional labor market with life cycle dynamics to account for income and risk inequality over the business cycle. Low-earners behave very differently from high-earners: initial wealth and search effort in the labor market have large implications for the large negative skewness of lifetime earnings for each of these two groups.

CO181 Room BCB 408 J-ISBA SESSION: NEW ADVANCES IN BAYESIAN STATISTICS

Chair: Beatrice Franzolini

C0240: **Bayesian Kolmogorov-Arnold neural model**

Presenter: **Myung Won Lee**, University of Edinburgh, United Kingdom

Co-authors: Miguel de Carvalho, Brian Reich

The aim is to present a novel statistical framework for Kolmogorov-Arnold neural networks (KAN), redefining them as a neural extension of generalized additive models. The canonical version of the model employs a doubly-additive approach in a three-layer configuration, while the deep version extends this to a deep additive model. The approach then motivates a novel picking-and-pruning prior, specifically designed for KAN architectures. This prior facilitates group-wise regularization of the spline coefficients governing the learnable activation functions, applied simultaneously to the inner and outer layers. Thus, the method not only performs variable selection but also streamlines the dense network, leveraging Bayesian inference to enhance interpretability for both prediction and classification tasks. The proposed method is validated with simulation studies on artificial data. The applicability of the model is further investigated on the chemical exposure and inflammation data.

C0297: **Diffusion piecewise exponential models for survival extrapolation using piecewise deterministic Monte Carlo**

Presenter: **Luke Hardcastle**, University College London, United Kingdom

Co-authors: Samuel Livingstone, Gianluca Baio

The purpose is to introduce the diffusion piecewise exponential model. Piecewise exponential models are a flexible non-parametric approach for time-to-event data, but extrapolation beyond final observation times typically relies on random walk priors and deterministic knot locations, resulting in unrealistic long-term hazards. The diffusion piecewise exponential model is a prior framework that consists of a discretized diffusion for the hazard and can encode a wide variety of information about the long-term behavior of the hazard, with time changed by a Poisson process prior for knot locations. This allows the behavior of the hazard in the observation period to be combined with prior information to inform extrapolations. Efficient posterior sampling is achieved using piecewise deterministic Markov processes, whereby we extend existing approaches using sticky dynamics from sampling spike-and-slab distributions to more general trans-dimensional posteriors. The focus is on applications in health-technology assessment, where the need to compute mean survival requires hazard functions to be extrapolated beyond the observation period, showcasing performance on datasets for Colon cancer patients.

C0451: **A unifying framework for generalized Bayesian online learning in non-stationary environments**

Presenter: **Gerardo Duran-Martin**, University of Oxford, United Kingdom

Co-authors: Leandro Sanchez-Betancourt, Alexander Shestopaloff, Kevin Murphy

A unifying framework is proposed for methods that perform probabilistic online learning in non-stationary environments. The framework is called BONE, which stands for generalized (B)ayesian (O)nline learning in (N)on-stationary (E)nvironments. BONE provides a common structure to tackle a variety of problems, including online continual learning, prequential forecasting, and contextual bandits. The framework requires specifying three modeling choices: (i) a model for measurements (e.g., a neural network), (ii) an auxiliary process to model non-stationarity (e.g., the time since the last changepoint), and (iii) a conditional prior over model parameters (e.g., a multivariate Gaussian). The framework also requires two algorithmic choices, which are used to carry out approximate inference under this framework: (i) an algorithm to estimate beliefs (posterior distribution) about the model parameters given the auxiliary variable, and (ii) an algorithm to estimate beliefs about the auxiliary variable. It is shown how the modularity of the framework allows for many existing methods to be reinterpreted as instances of BONE, and it allows the proposal of new methods.

CC420 Room BCB 208 ESG SCORES AND SUSTAINABLE FINANCE

Chair: Massimiliano Caporin

C0289: **ESG controversies and market recovery: A wavelet perspective**

Presenter: **Michaela Kiermeier**, University of Applied Sciences Darmstadt, Germany

Capital-market reactions to ESG controversies can be explained by various models, e.g., signaling, institutional legitimacy, and capital market inefficiencies. Together, these perspectives suggest that scandals around environmental, social, or governance may trigger short-run underperformance that can revert into longer-term recovery once effective countermeasures are perceived. To test this empirically, wavelet analysis is applied to trace how information on ESG controversies affects the financial performance of investment strategies for publicly listed European firms. Event windows are defined by controversial signals from the news. Expected returns are modeled using the best-fitting multifactor risk models for each company; abnormal returns are then computed around each event. Maximal overlap discrete wavelet transforms decompose these abnormal returns across multiple time scales, allowing the identification of transient versus persistent effects. A bootstrap procedure generates scale-specific thresholds so that only statistically significant wavelet coefficients are retained, ensuring robust inference about the duration and magnitude of ESG impacts. By

combining a theoretically grounded view of controversy resolution with a scale-by-scale time-frequency approach, it is illuminated whether, when, and for how long ESG information is priced, offering new insights into the temporal transmission mechanisms between corporate responsibility shocks and asset prices.

C0344: Reconciling returns data and ESG criteria: Integrating responsibility into mean-variance portfolio construction

Presenter: **Tomer Shushi**, Ben Gurion University of the Negev, Israel

Socially responsible portfolio selection has recently gained much attention. Environmental, social, and governance (ESG) scores offer a means to quantify a company's contributions to both the environment and society. Novel frameworks are proposed for optimal portfolio selection that provide a tradeoff between entirely focusing on the historical data of each of the stocks in the portfolio and entirely focusing on their ESG scores. To obtain the optimal weights as part of the portfolio selection process, a multivariate constrained optimization problem needs to be solved. An explicit solution is proposed for such a problem, and it is explored using empirical illustrations for the mean-variance, tail-value-at-risk measures, and other classes of risk measures.

C0189: Retirement income phased withdrawal plans and the gender pension gap

Presenter: **Alexandra Dias**, University of York, United Kingdom

The current shift from guaranteed benefits to phased withdrawal pension plans moves the financial risk from the corporate sector to households, potentially contributing to the gender pension gap. Women living in retirement on a lower income than men is a well-known problem. The aim is to study how investment decisions, level of consumption, and longevity affect the welfare of men and women during retirement. It is found that a gender pension gap is inherent to phased withdrawal plans, and that even more aggressive, riskier investment strategies do not successfully close the gap. It is concluded that the currently popular phased withdrawal retirement plans contribute to perpetuating the gender pension gap.

CC433 Room BCB 307 HIGH-DIMENSIONAL AND SPARSE MODELING

Chair: Andreas Artemiou

C1149: Sparse estimation in semi-parametric finite mixture of regression models subject to right censoring

Presenter: **Farhad Shokoohi**, University of Nevada Las Vegas, United States

The aim is to propose a sparse estimation framework for semi-parametric finite mixture of regression models in the presence of right-censored outcomes. Such models are particularly relevant for heterogeneous populations where subgroups exhibit distinct regression relationships, yet the underlying error distributions are left unspecified. The approach integrates a penalized likelihood method with sieve-based nonparametric density estimation to accommodate both component heterogeneity and censoring. The sparsity-inducing penalty facilitates variable selection within each mixture component, enhancing interpretability while maintaining predictive accuracy. The asymptotic properties of the proposed estimator are established, including selection consistency and oracle efficiency, under mild regularity conditions. Extensive simulation studies demonstrate superior performance in identifying relevant covariates and recovering component structures compared to existing methods. An application to a real survival dataset illustrates the practical utility of the method in uncovering latent subpopulation-specific effects in censored regression settings.

C1258: An algorithm for solving the constrained sparse-group lasso

Presenter: **Nazgul Zakiyeva**, Institute of Mathematics and Mathematical Modelling, Kazakhstan

The aim is to propose a mixed coordinate descent method of multipliers for solving the constrained sparse-group lasso problem in multivariate time series analysis. The constrained sparse-group lasso extends the widely used sparse-group lasso by incorporating linear constraints, which enable the integration of prior information into the model, such as balance conditions between supply and demand. The proposed method combines coordinate descent with augmented Lagrangian techniques, yielding an efficient and flexible optimization framework that simultaneously addresses sparsity, group structure, and linear restrictions. The methodology is examined through numerical experiments and real data applications, and benchmarked against alternative optimization approaches.

C1394: Heuristic algorithms for subset regression model selection

Presenter: **Cristian Gatu**, Alexandru Ioan Cuza University of Iasi, Romania

Co-authors: Georgiana-Elena Pascaru, Petru Sebastian Drumia, Erricos Kontoghiorghe

Heuristics step-wise algorithms are an established approach to the regression model selection problem. A main drawback of these methods is the reduced number of submodels that are evaluated in order to select a solution, thus, in general, failing to find the optimum. Three strategies that aim to overcome this issue are investigated: This first one (SEL- k) builds on standard forward selection, but selects at each step the best k variables, instead of the only best one. The second method (TREE- k) explores a combinatorial search space, thus increasing the number of submodels that are investigated. Specifically, at each step of the algorithm a new search branch is considered for each of the best k most significant variables. A branch will terminate either when there are no more significant variables to choose from, or when all variables have been considered. The third one (SHUTTLE- k) aims to obtain good solutions but avoiding a prohibitive computational cost. A list of k submodels is stored. During one iteration, for each of the k submodels, the best k variables are chosen yielding k^2 submodels (augmentation step). From the resulting k^2 list the best k submodels are kept for the subsequent iteration (reduction step). Various experiments are conducted on both real and artificially generated datasets in order to assess the three proposed algorithms. The results are presented and discussed.

CC407 Room BCB 313 EMPIRICAL FINANCE

Chair: Jonas Andersson

C0620: Skewed interest rate expectations and effects of central banks' market operations: Evidence from granular data

Presenter: **Takatashi Sasaki**, Bank of Japan, Japan

Co-authors: Taihei Sone, Daisuke Miyakawa, Kohei Maehashi

Using trade repository data on transaction records of Japanese yen-denominated overnight index swaps, individual market participants' expectations on future interest rates are estimated, and their time-variant distribution is documented with its higher order moments. By leveraging this novel information, quantitative exercises are implemented to verify the state-dependent effects of the Bank of Japan (BoJ)'s outright purchase of Japanese Government Bonds (JGBs) on the JGB yields conditional on the moments of this expectation distribution. It is found that the BoJ's fixed-rate purchase operation resulted in a larger reduction of the JGB yields when the expectation distribution on future interest rates was skewed more positively. This empirical result implies the usefulness of the estimated expectation distribution for central banks to conduct market operations effectively.

C1207: Reconsidering success in staged venture investments - The role of follow-on investment in investor lifecycles

Presenter: **Eugene Hong**, Keio University, Japan

Co-authors: Hiroshi Takahashi

The purpose is to suggest a different perspective on the success of staged venture investments for investors. Previous studies have primarily considered exit events, such as IPOs and M&A, as the measure of success in venture investments. On the other hand, the success of venture investment should not be limited to the event of converting invested capital into multiplied cash as a one-time occurrence. Throughout the investor's lifetime of investing activities, it is considered that the firm's successful serial fundraising can be an essential event that determines the sustained success of the firm's investors. From this perspective, it is examined how an investor's participation in a follow-on investment in the portfolio company affects the likelihood of their long-term success across the lifetime investment activities. Furthermore, methods for predicting the macroscopic success

of venture investments are proposed using machine learning techniques, trained on both the investors' past investment history data and the firms' fundraising records.

C1332: Network risk spillovers between European energy related firms

Presenter: **Jone Ascorbebeitia**, University of the Basque Country UPV/EHU, Spain

Co-authors: Jordi Barbera, Susan Orbe

The purpose is to investigate risk spillovers among Eurozone energy-related industries over the period 2015-2025. The energy transition, combined with geopolitical and economic shocks, is reshaping European stock markets, altering sector dynamics and amplifying sectoral vulnerabilities. Understanding how extreme events propagate across the energy chain is therefore crucial for both investors and regulators. The focus is on the most capitalized firms in the Energy, Utilities, and Basic Materials industries, which represent the production, distribution, and consumption stages of energy in the Eurozone. Using the tail event driven network (TENET) framework, the transmission of tail risks is captured, and systemic firms both within and across industries are identified. Findings reveal three pronounced waves of spillovers, corresponding to the trade conflict between U.S. and China, the COVID-19 outbreak, and the Russia-Ukraine war. The results show that the upstream Energy industry consistently acts as a net transmitter of risk, while Utilities and Basic Materials mainly absorb shocks. Moreover, spillovers are primarily driven by intra-industry contagion, underscoring the importance of sector-specific connections. Overall, the analysis highlights that systemic vulnerabilities extend beyond financial institutions and that energy-related industries play a critical role in the stability of Eurozone markets during the ongoing energy transition.

CC431 Room BCB 402 APPLIED FINANCIAL ECONOMETRICS

Chair: Johan Lyhagen

C0445: From imports to fracking: A Markov-switching Granger causality analysis of crude oil prices and political polarization

Presenter: **Charisios Grivas**, Aalborg University, Denmark

Co-authors: Rubens Morita

Political polarization has emerged as a major source of macroeconomic uncertainty, yet its connection to commodity markets is not well understood. The evolving relationship between U.S. partisan conflict and crude oil returns is examined through a Markov-switching Granger causality framework with time-varying volatility and correlation, which permits the endogenous identification of structural breaks in the predictive linkages. Based on monthly data from 1981-2025, the model identifies two distinct regimes: Before 2009, oil shocks systematically Granger-caused increases in partisan conflict, reflecting the destabilizing role of import dependence and global supply disruptions; after the global financial crisis and the shale boom, this predictive channel collapsed, and oil prices no longer convey systematic information about domestic polarization. These results challenge the conventional wisdom that oil prices serve as a barometer of U.S political conflict, showing instead that their informational content is contingent on the broader energy and political environment.

C1238: Safe distance to systemic risk

Presenter: **Sylvain Benoit**, University Paris Dauphine - PSL, France

Co-authors: Renzhi Liu

The aim is to propose a new systemic risk indicator to measure the distance to the extreme losses of a financial system. Constructed from daily out-of-sample value-at-risk (VaR) exceptions across 95 large U.S. financial institutions from 2000 to 2023, the indicator calculates the shortfall in market value during these exceptions. By applying extreme value theory (EVT) to the maximum weekly shortfalls using a half-year rolling window, we effectively model the tail risk of the financial system. The empirical analysis demonstrates that this indicator captures significant financial crises accurately, such as the Great Financial Crisis of 2008, the sovereign debt crisis of 2010, and the COVID-19 pandemic in 2020. Through quantile regression, it is shown that increases in the indicator significantly predict negative shocks to industrial production growth rates.

C1280: Seasonal and regime-dependent effects of heat stress on European tourism equities

Presenter: **Iason Kyrlis**, University of Piraeus, Greece

Co-authors: Charalampos Agiropoulos, George Galanos, Sotirios Varelak, Christos Kanellopoulos

The impact of climate-induced heat stress is investigated, measured by the Universal Thermal Climate Index (UTCI), on European tourism equities. Using weekly data from 2010-2025, Newey-West OLS regressions are combined with time-frequency wavelet methods. Econometric results show that anomalous-day frequencies in the Northeast quadrant significantly predict tourism portfolio returns (+22 basis points per day), while summer (JJA) heat is associated with negative performance (-39 bp/). A structural shift emerges post-2020, with positive sensitivities (+46 bp/). Wavelet coherence analysis with AR(1) Monte Carlo significance reveals statistically reliable short-horizon (2-8 week) co-movements, with returns typically leading heat stress by 0.5 weeks, consistent with anticipatory pricing. However, coherence is sparse in summer-only samples, suggesting that the seasonal drag reflects broad demand effects rather than concentrated bursts of extreme heat. These findings highlight spatial heterogeneity, regime dependence, and the importance of multi-resolution methods in quantifying climate-financial linkages.

CC436 Room BCB 405 APPLIED STATISTICAL TOOLS AND EVALUATION

Chair: Steven Gilmour

C0308: A flexible framework for adaptive designs based on the simulated annealing algorithm

Presenter: **Francesco Mariani**, University of Bologna, Italy

Co-authors: Rosamarie Frieri, Marco Novelli

In adaptive randomization designs, new trial participants are sequentially assigned to a particular treatment based on observed covariates and/or responses. While accounting for covariate profiles can increase statistical power and efficiency, solely relying on them may cause ethical concerns since patients are potentially allocated to the inferior treatment. A sequential procedure is proposed based on the simulated annealing (SA) algorithm that allows combining different research needs by allocating patients in a covariate-adjusted (CA) and/or response-adaptive (RA) manner. Thanks to SA flexibility, the resulting procedure can be: (i) entirely RA, when patients are assigned to the best-performing treatment; (ii) RA + CA, where allocation also aims at balancing covariates across treatments; (iii) CARA, assigning each patient to the most effective treatment based on their covariate profile, thus reflecting a personalized medicine approach. The procedure can be either model-free or model-based, so as to work at a population or patient-specific level, and thus depending on the assumptions over treatment-covariate interactions. In complex scenarios where treatment effect varies with covariates, Bayesian additive regression trees (BART) are incorporated to guide SA toward optimal allocations. The resulting class of CARA designs has appealing statistical and ethical properties, as illustrated through simulation studies.

C1372: On statistical reproducibility of ROC-based diagnostic tests

Presenter: **Hamdah Alshamari**, Durham University, United Kingdom

Co-authors: Tahani Coolen-Maturi, Frank Coolen

Hypothesis testing based on diagnostic measures is widely applied in medical science and healthcare, yet repeated testing can lead to different conclusions. Assessing the statistical reproducibility of such tests is therefore essential to ensure reliable results. Reproducibility probability (RP) measures the probability that the same outcome would be obtained if a test were repeated under identical conditions with the same sample size. Nonparametric predictive inference (NPI) offers a predictive framework for studying RP, using the NPI bootstrap method to evaluate RP for accuracy tests. Unlike traditional methods, NPI focuses on prediction rather than estimation within a frequentist framework. RP is applied to the area under the ROC curve (AUC) test, with a simulation study illustrating the NPI reproducibility approach. The AUC is a widely used metric for evaluating the performance of diagnostic tests in distinguishing between diseased and non-diseased individuals. The findings show that the

RP of AUC tests can be low, particularly when the p-value is near the significance threshold, raising important concerns about reproducibility in diagnostic test evaluation.

C0210: Statistical tools for assessing mixtures effects on health outcomes: A trade-off between complexity and interpretability

Presenter: **Susana Diaz Coto**, Geisel School of Medicine at Dartmouth, Dartmouth Health, United States

Several studies have explored the joint effect of mixtures of metals on health outcomes by using sophisticated statistical techniques, such as Bayesian kernel machine regression (BKMR). Although this method can detect complex relationships between metals, final conclusions are usually simplified in terms of an increase/decrease of the risk function at fixed values for individual metals. The feasibility and interpretability of BKMR are analytically and graphically explored to assess the effect of metal mixture exposures on neuropsychological development in early childhood. The scores derived by the BKMR exposure-response function were computed and compared with those provided by traditional linear regression models (LRM). Pearson correlation coefficients of the estimations obtained by BKMR and LRM accounting for pairwise interactions between metals ranged from 0.92 to 0.95, for the three previously identified latent domains: Executive functions, motor functions, and visual and verbal functions. The observed differences between both estimations mainly occurred in participants having the lowest or highest values of individual metals, which suggests linearity in the associations and not high-order interactions. It is concluded that, in this case, employing linear regression to model the impact of mixtures on targeted outcomes led to similar results to more complex techniques, allowing a better understanding of both individual and interaction effects between metals.

CC408 Room BCB 406 METHODOLOGICAL STATISTICS

Chair: Sanjoy Sinha

C1111: Efficient estimation for scale parameter of Birnbaum Saunders distribution for multiple dimensions

Presenter: **Waqas Makhdoom**, Department of Statistics, Government College University, Lahore, Pakistan

Co-authors: Muhammad Kashif Ali Shah, Nighat Zahra, Ejaz Ahmed

The Birnbaum-Saunders distribution is one of the most reputable probability distributions for fatigue life and reliability studies, having shape and scale parameters. The parameter estimation of this distribution is a common area of interest for researchers in recent times. The maximum likelihood estimation approach is generally used for point estimation due to its appealing asymptotic properties. The interest is in boosting the efficiency of the scale estimator while incorporating non-sample information. Some improved estimation strategies are employed, such as the restricted linear shrinkage estimator, the preliminary test estimator, the shrinkage preliminary test estimator, and the Stein-type shrinkage estimators. The asymptotic properties of the suggested estimators are derived both in analytical and numerical terms. The graphical presentation of the performance of the estimators is also presented. A real-life data set is analyzed to judge the performance of the proposed estimators.

C1379: Laplace approximations for Gaussian process and mixed effects quantile regression

Presenter: **Andrea Thomas Nava**, Hochschule Luzern, Switzerland

Co-authors: Fabio Sigrist

The aim is to propose novel Laplace approximations for quantile regression with Gaussian process and mixed effects models that address limitations of standard Laplace approximations by replacing the Hessian with data-adaptive approximations. For large-scale applications with Gaussian processes, Vecchia approximations are used to enable scalable inference. The quality of the methods concerning posterior and marginal likelihood approximations is evaluated on both simulated and real-world datasets, benchmarking against existing state-of-the-art approaches for quantile regression with Gaussian processes and grouped random effect models. The methods are found to be computationally more efficient and stable while often providing more accurate posterior predictive distributions and hyperparameter estimates.

C1414: Scalable covariance parameter estimation of nonstationary Matern process from high-resolution spatial data

Presenter: **Kunal Das**, Iowa State University, United States

Co-authors: Zhengyuan Zhu

Accurately modeling nonstationary spatial dependence is crucial for high-resolution spatial data from satellite imaging and sensors. The aim is to tackle these challenges by estimating key parameters of nonstationary Gaussian random fields under infill asymptotics, explicitly focusing on the smoothness parameter and spatially varying microergodic functions within a Matern covariance framework. A divide-and-conquer strategy has been proposed utilizing recent advancements in higher-order quadratic variations with multivariate discrete differentiation. By overlaying a coarse grid of anchor points on the observation domain, locally stationary neighborhoods are constructed to estimate unknown but constant stationary parameters. Then, the spatially varying structure of the concerned parameters is obtained using kernel smoothing over the domain. The methodology provides theoretical guarantees for the asymptotic consistency of local and global estimators, without requiring prior knowledge of model parameters and under mild smoothness conditions of unknown component functions. This framework is scalable, interpretable, statistically robust, and relevant for geostatistical modeling applications. A two-fold simulation study is also discussed to assess the finite sample accuracy of asymptotic theoretical results and to illustrate the computational efficiency of the method compared to existing techniques for large spatial datasets.

CC429 Room BCB 409 HIGH-FREQUENCY AND STOCHASTIC ECONOMETRICS

Chair: Masayuki Uchida

C0796: Test for volatility parameter change in linear parabolic SPDEs

Presenter: **Yozo Tonaki**, The University of Osaka, Japan

Co-authors: Yusuke Kaino, Masayuki Uchida

The purpose is to consider change point detection for the volatility parameter in second-order linear parabolic stochastic partial differential equations (SPDEs) based on high-frequency spatiotemporal data. Ornstein-Uhlenbeck processes are obtained from the inner product of the solution of the SPDE and an orthonormal basis, and these processes can be approximated through statistical inference for SPDEs based on high-frequency spatiotemporal data. A test statistic is therefore proposed to detect changes of the volatility parameter in the linear parabolic SPDE based on change point analysis for diffusion processes. The asymptotic null distribution of the proposed test statistic is derived, and the test statistic is shown to be consistent. Additionally, the asymptotic properties of the test statistic are validated through numerical simulations.

C1126: QBIC of SEM for diffusion processes based on high-frequency data

Presenter: **Shogo Kusano**, Kumamoto University, Japan

Co-authors: Masayuki Uchida

Structural equation modeling (SEM) is a statistical method to investigate relationships among latent variables. Since SEM is a confirmatory analysis method, the model needs to be specified in advance based on the theoretical framework of the respective research field. However, statisticians may often have several candidate models for SEM and must choose the optimal one from among them. To address this issue, various information criteria for SEM have been actively studied. In particular, the Akaike information criterion (AIC) and the Bayesian information criterion (BIC) are commonly used. Recently, SEM for diffusion processes based on high-frequency data has been studied. For model selection of the SEM, an AIC-type statistic based on the quasi-likelihood has been proposed. However, this criterion does not ensure model selection consistency. Therefore, BIC-type statistics are considered for SEM. The asymptotic expansion of the marginal quasi-log likelihood is first obtained. Based on this result, two types of quasi-BIC are proposed for SEM, and it is shown that the information criteria have model selection consistency. Furthermore, some examples are provided, and simulation studies are conducted.

C1226: Asymptotically uniformly most powerful tests for diffusion processes with nonsynchronous observations

Presenter: **Teppei Ogihara**, University of Tokyo, Japan

The purpose is to introduce a quasi-likelihood ratio testing procedure for diffusion processes observed under nonsynchronous sampling schemes. High-frequency data, particularly in financial econometrics, are often recorded at irregular time points, challenging conventional synchronous methods for parameter estimation and hypothesis testing. To address these challenges, a quasi-likelihood framework is developed that accommodates irregular sampling while integrating adaptive estimation techniques for both drift and diffusion coefficients, thereby enhancing optimization stability and reducing computational burden. The asymptotic properties of the proposed test statistic are rigorously derived, showing that it converges to a chi-squared distribution under the null hypothesis and exhibits consistency under alternatives. Moreover, it is established that the resulting tests are asymptotically uniformly most powerful. Extensive numerical experiments corroborate the theoretical findings and demonstrate that the method outperforms existing nonparametric approaches.

Sunday 14.12.2025

13:40 - 15:20

Parallel Session H – CFE-CMStatistics 2025

CI010 Room BCB 211 DYNAMIC FACTOR MODELS AND BEYOND**Chair: Esther Ruiz****C0173: Identification and estimation of dynamic factor models****Presenter:** Joerg Breitung, University of Cologne, Germany**Co-authors:** Matei Demetrescu

In empirical practice, dynamic factor models are typically estimated by performing a principal component analysis (PCA) on the static factor representation. A two-step PCA estimator and a sequential least-squares approach are proposed to estimate the original dynamic factors that enter the model with a prespecified number of lags. The identification of the dynamic factors subject to the usual normalization restrictions is analyzed, and the consistency of the estimator is established. Furthermore, consistent methods for selecting the number of dynamic factors and the number of lags are discussed.

C0174: Forecasting ENSO: A frequency band factor model approach**Presenter:** Alessandro Giovannelli, University of L'Aquila, Italy**Co-authors:** Tommaso Proietti

The El Nino-Southern Oscillation (ENSO) is a major driver of interannual climate variability, with significant impacts on global climate patterns. The aim is to introduce a novel forecasting methodology based on dynamic factor models, where the unobserved factors are estimated separately for two distinct frequency bands: One capturing ENSO-related medium-to-long-term fluctuations (17 years), and another isolating short-term variability. The rationale for this frequency decomposition is that separating signals by frequency enhances medium-term predictive information, while still retaining potentially useful information from short-term fluctuations. Once the factors are estimated, independently for each frequency band, they are incorporated into a sparse regression forecasting model, which dynamically selects the most predictive components. An extensive real-time forecasting application using over 2000 climate variables illustrates the potential benefits of this frequency-specific approach for ENSO prediction.

C0175: IRF and nowcasting with functional approaches and mixed-frequency data**Presenter:** Catherine Doz, Paris School of Economics, France**Co-authors:** Laurent Ferrara, Anna Simoni

Prediction and causal models often involve economic time series with different sampling frequencies. The mix-frequency problem has received a lot of attention in the past nowcasting/forecasting literature, which has proposed solutions that work especially well for prediction when the frequency gap is small, like the gap between monthly and quarterly data. An alternative approach is proposed to deal with mix-frequency which is particularly designed for situations where the gap in sampling frequencies is large, like daily and quarterly. Such a large gap can be easily found when one mixes series from standard and non-standard sources, like the internet or newspapers. The approach focuses on both prediction and causality, and exhibits excellent performance. By treating the high-frequency variable as a realization of a stochastic process in continuous time, the estimation problem is cast in the class of ill-posed inverse problems, and different regularized estimators are proposed. Importantly, for each high-frequency covariate, a measure of its influence is recovered on the target variable as a function of the time gap between the forecasting horizon and the date of the information in the past. Because the analysis is conditional on several covariates, the fact that the influence of high-frequency series rapidly decreases as soon as information from standard data arises is well captured.

CO068 Room BCB G07 HiTEC: ADVANCES IN FINANCIAL ECONOMETRICS**Chair: Genaro Sucarrat****C0810: Noncausal AR processes driven by causal GARCH volatility****Presenter:** Jean-Michel Zakoian, CREST, France**Co-authors:** Daniel Velasquez-Gaviria

The purpose is to study the introduction of causal conditional heteroskedasticity in noncausal autoregressive (AR) models. It is demonstrated that large shocks to the independent innovation that drives the GARCH error term of a noncausal first-order AR model result in heightened volatility following a bubble crash. The non-coincidence of the information sets generated by past observations and past values of the GARCH process makes estimation non-standard. In particular, the full quasi-maximum likelihood estimator (QMLE) is generally inconsistent. The asymptotic properties of three-step weighted least squares estimators of the AR coefficient and the QMLE of the volatility coefficients are investigated. Findings are illustrated via Monte Carlo experiments and real financial data.

C0896: An economic evaluation of exchange rates higher order moments timing**Presenter:** Francesco Violante, IESEG School of Management, France**Co-authors:** Stefano Grassi, Mattia Alfiero

The sequential Monte Carlo and its online variant are proposed for the estimation of multivariate Garch models that feature time-varying skewness and kurtosis. Relative to model-specific Markov chain Monte Carlo, sequential Monte Carlo has the advantages of generality, parallelizability, and speed. Moreover, it gives as output the marginal likelihood useful in model selection. In the empirical application, the forecastability of the exchange rate is revisited through an asset allocation problem, using different specification of multivariate volatility models.

C0807: Climate-driven shockwaves along the futures curve**Presenter:** Susana Campos Martins, Catholic University of Portugal, Portugal**Co-authors:** Filippo Pellegrino, Jose Afonso Faia

Climate-driven disasters are increasingly disrupting global commodity supply chains, yet the transmission of these shocks beyond the spot market remains poorly understood. The aim is to pioneer a granular, county-level geospatial dataset that maps production sites, storage facilities, and transport hubs for key energy commodities from 1970 to 2025. By aligning these locations with high-resolution records of major climate and weather disasters, long-standing data constraints are overcome that have sidelined commodities from mainstream climate finance. Building on the theory of storage, it is quantified how climate-driven disasters jointly affect spot prices and the convenience yield, revealing under what conditions long-dated futures rise above pre-shock levels following a disruption.

C0792: Prediction of non-negative variables under misspecification and nonstationarity**Presenter:** Genaro Sucarrat, BI Norwegian Business School, Norway**Co-authors:** Christian Francq

In empirical practice, it is commonly assumed that the entertained model is equal to the conditional expectation. While this simplifies theory and interpretation, it is unlikely to hold in reality. The purpose is to consider a broad class of predictive specifications for non-negative variables (e.g., volatility, spreads, volume, and unemployment). The predictive specifications need not equal the (unknown) conditional expectation, and the non-negative variables can be nonstationary. Examples of predictive specifications in the class include ARMA specifications and nonstationary seasonality predictors (and combinations thereof). Consistent and asymptotically normal estimation of the parameters of the predictive specifications is established for the quasi maximum likelihood estimator (QMLE). The finite sample properties of the estimator is studied by simulation, and an empirical application illustrates the results.

CO059 Room BCB G08 APPLIED MACRO**Chair: Michael Owyang****C0227: Identifying business cycles in real time with quantile dynamic factor Markov switching models***Presenter:* **Jeremy Piger**, University of Oregon, United States

The global pandemic created extreme swings in macroeconomic data that complicate statistical analysis. Simply dummifying out these observations is not a good option, as this would eliminate macroeconomic objects of interest, such as recessions. The existing literature has focused on model augmentation in the form of time-varying volatility (e.g., stochastic volatility, ARCH, GARCH, etc.) Such approaches have been taken to workhorse models in empirical macroeconomics, such as VARs and dynamic factor models. An alternative approach is evaluated for Markov-switching models, which are a popular class of models designed to identify recessions in economic data. Specifically, the performance of dynamic factor Markov switching (DFMS) models is investigated, which are specified in terms of conditional quantiles of recessions rather than conditional means. It is shown that these quantile DFMS models are capable of capturing recessions, both retrospectively and in real time, while not being too strongly influenced by extreme recessions and expansions, such as those observed during and following the March-April 2020 recession.

C0453: Asymmetric loss and financial market forecasts*Presenter:* **Julie Bennett**, Duke University, United States*Co-authors:* Michael Owyang

The relationship between inflation and interest rate forecasts made by financial markets is evaluated. Using forecast rationality testing under univariate and multivariate asymmetric loss functions, evidence is found that financial market forecasts of inflation are made conditional on the expected path of monetary policy. Further, it is found that the degree of asymmetry financial markets exhibit depends on the current level of the inflation rate. These results provide context for how financial market forecasts should be interpreted and incorporated into larger economic models.

C0738: The macroeconomic impact of unexpected data releases*Presenter:* **Ana Beatriz Galvao**, Bloomberg Economics, United Kingdom

Key macroeconomic data announcements demand abnormal compensation for risk. The aim is to investigate whether the market reaction to data surprises can be explained by their contribution to improving our understanding of the current and future state of the economy. It also examines whether the market reaction plays a role in the transmission of structural macroeconomic shocks.

C0578: Does uncertainty forecast recessions?*Presenter:* **Michael Owyang**, Federal Reserve Bank of St Louis, United States

Increases in uncertainty tend to be highly correlated with declines in current and future economic activity. The literature on the real-time predictive ability of uncertainty for economic activity is smaller, but the findings are similar. It is not surprising, then, that incorporating policy uncertainty into binary outcome models has been argued to improve recession forecasts. A prior study evaluates whether the economic policy uncertainty index (EPU) and its components have predictive ability over other indicators such as the term spread and the growth rate of the S&P500. Another study uses the uncertainty index from a past study, performing their experiments in-sample. Alternative measures of uncertainty are evaluated as a predictor of business cycles using quasi-real-time experiments. A real-time estimate of the prior study's uncertainty measure is constructed for comparison with the real-time measure of another study. It is found that a combination measure of the two indices calibrated for business cycle prediction best forecasts future recessions.

CO058 Room BCB G09 STATISTICAL AND ECONOMETRIC MODELS FOR CENSORED DATA**Chair: Ralf Wilke****C0847: Censored lifespans in a double-truncated sample: Maximum likelihood inference for exponential and geometric***Presenter:* **Fiete Sieg**, Universität Rostock, Germany*Co-authors:* Rafael Weissbach, Anne-Marie Toparkus

A small truncated sample of censored durations is approximated with a bivariate Poisson process. For exponential or geometric distribution of the duration, the likelihood is derived, profile out the unobservable sample size, and study identification, consistency, as well as asymptotic normality.

C0622: Multiple contrast tests for RMTLs in factorial competing risks settings*Presenter:* **Merle Munko**, Otto-von-Guericke University Magdeburg, Germany*Co-authors:* Dennis Dobler, Marc Ditzhaus

Easy-to-interpret effect estimands are highly desirable in survival analysis. In the competing-risks framework, one good candidate is the restricted mean time lost (RMTL). It is defined as the area under the cumulative incidence function (CIF) up to a prespecified time point and, thus, it summarizes the CIF into a meaningful estimand. While existing RMTL-based tests are limited to two-sample comparisons and two risks, we aim to develop general contrast tests for factorial designs and an arbitrary number of risks based on a Wald-type test statistic. Furthermore, the often-made, rather restrictive continuity assumption is avoided on the event time distribution. This allows for ties in the data, which often occur in practical applications, e.g., when event times are measured in (whole) days. In addition, more reliable tests are developed for RMTL comparisons that are based on resampling procedures to improve the small sample performance. In a second step, multiple tests for RMTL comparisons are developed to test several null hypotheses simultaneously. The asymptotically exact dependence structure is incorporated between the local test statistics to gain more power.

C0719: Tests of exogeneity in duration models with censored data*Presenter:* **Gilles Crommen**, KU Leuven, Belgium*Co-authors:* Ingrid Van Keilegom, Jean-Pierre Florens

Consider a duration time of interest T , a possibly endogenous treatment variable Z , and a vector of exogenous covariates X such that $T = \varphi(Z, X, U)$ is increasing in U with $U \sim U[0, 1]$. Moreover, let T be right-censored by a censoring time C such that only their minimum, denoted by the follow-up time $Y = \min\{T, C\}$, is observed. Test statistics are constructed for the hypothesis that Z is exogenous w.r.t. T , that is, Z is independent of U . The tests make use of an instrumental variable W that is independent of U given X , since it can be shown that Z is exogenous w.r.t. T if and only if $V_T = F_{T|Z,X}(T | Z, X)$ is independent of (W, X) jointly. The asymptotic properties of the proposed test are proved for the case where Z, X , and W are categorical, and possible bootstrap approximations for the critical value of the tests are shown to have a good finite sample performance via simulations. Lastly, an empirical example is provided using data from the National Job Training Partnership Act (JTPA) Study.

C0720: On the identifiability under dependent censoring or in a competing risk setting using divergence measures*Presenter:* **Roel Braekers**, Hasselt University, Belgium

In a competing risk setting or under dependent censoring, latent random variables are often assumed to model the different times until certain failure events. However, this type of data is observed only for the minimum event time and the cause of failure. Therefore, it needed to make assumptions on the association between the latent random variables to identify the joint survival function of the latent random variables from the observed quantities. For example, assuming independence or a known association function (in the form of a copula function) between these random variables is a common assumption in a competing risk setting to avoid identifiability. When a parametric family of copula functions is assumed for the association and parametric marginal distributions for the different event times, it is harder to check whether the assumed model is identifiable from the observed quantities. Divergence measures are used to develop a method to verify this identifiability in a competing risk setting or under dependent censoring. It is shown that this method works for various parametric copula functions for the association between the different event

times and their parametric marginal distributions, even when the number of parameters increases to a very large number.

CO304 Room Virtual R01 RECENT ADVANCES IN CAUSAL INFERENCE
Chair: Soham Jana
C0235: Distilling heterogeneous treatment effects: Stable subgroup estimation in causal inference

Presenter: **Tiffany Tang**, University of Notre Dame, United States

Recent methodological developments have introduced new black-box approaches to better estimate heterogeneous treatment effects; However, these methods fall short of providing interpretable characterizations of the underlying individuals who may be most at risk or benefit most from receiving the treatment, thereby limiting their practical utility. Causal distillation trees (CDT) are introduced to estimate interpretable subgroups. CDT allows researchers to fit any machine learning model to estimate the individual-level treatment effect, and then leverages a simple, second-stage tree-based model to "distill" the estimated treatment effect into meaningful subgroups. As a result, CDT inherits the improvements in predictive performance from black-box machine learning models while preserving the interpretability of a simple decision tree. Theoretical guarantees are derived for the consistency of the estimated subgroups using CDT, and stability-driven diagnostics are introduced for researchers to evaluate the quality of the estimated subgroups. The proposed method is illustrated on a randomized controlled trial of antiretroviral treatment for HIV from the AIDS Clinical Trials Group Study 175, and it is shown that CDT outperforms state-of-the-art approaches in constructing stable, clinically relevant subgroups.

C0250: Fisher consistency of surrogate losses for optimal dynamic treatment regimes

Presenter: **Nilanjana Laha**, Texas A&M University, United States

Co-authors: Nilson Chapagain, Aaron Sonabend

Patients with chronic diseases often receive treatments over multiple stages. The aim is to learn the optimal dynamic treatment regime (DTR) from longitudinal data, where the number of stages and treatment options per stage are arbitrary. This reduces to a sequential, weighted multiclass classification problem. This is addressed by solving the classification problem across all stages using Fisher consistent surrogate losses. While special cases (e.g., binary treatments) admit such surrogates, a general theory remains undeveloped. Necessary and sufficient conditions are established for DTR Fisher consistency within the class of non-negative, stagewise separable surrogate losses, offering the first such result in the DTR literature. It is further shown that many convex surrogates, including smooth, permutation-equivariant, and relative-margin-based ones, are inconsistent in this setting. To overcome this, SDSS (simultaneous direct search with surrogates) is proposed, which leverages smooth, non-concave surrogates to learn optimal DTRs. A gradient-based algorithm is introduced for SDSS, and a sharp regret bound is derived under a small optimization error.

C0591: A debiased estimator for the mediation functional in (ultra) high-dimension in the presence of interaction effects

Presenter: **Debarghya Mukherjee**, Boston University, United States

Co-authors: AmirEmad Ghassami, Shi Bo

Mediation analysis is a crucial tool for uncovering the mechanisms through which a treatment affects the outcome, providing deeper causal insights and guiding effective interventions. Despite advances in analyzing the mediation effect with fixed/low-dimensional mediators and covariates, the understanding of estimation and inference of mediation functional in the presence of (ultra)-high-dimensional mediators and covariates is still limited. An estimator is presented for mediation functional in a high-dimensional setting that accommodates the interaction between covariates and treatment in generating mediators, as well as interactions between both covariates and treatment and mediators and treatment in generating the response. It is demonstrated that the estimator is \sqrt{n} -consistent and asymptotically normal, thus enabling reliable inference on direct and indirect treatment effects with asymptotically valid confidence intervals. A key technical contribution is the development of a multi-step debiasing technique, which may also be valuable in other statistical settings with similar structural complexities where accurate estimation depends on debiasing. The proposed methodology is evaluated through extensive simulation studies and is applied to the TCGA lung cancer dataset to estimate the effect of smoking, mediated by DNA methylation, on the survival time of lung cancer patients.

C1022: Stochastic optimization algorithms for instrumental variable regression with streaming data

Presenter: **Abhishek Roy**, Texas A&M University, United States

Co-authors: Xuxing Chen, Yifan Hu, Krishnakumar Balasubramanian

The aim is to develop and analyze algorithms for instrumental variable regression by viewing the problem as a conditional stochastic optimization problem. In the context of least-squares instrumental variable regression, the algorithms neither require matrix inversions nor mini-batches, thereby providing a fully online approach for performing instrumental variable regression with streaming data. When the true model is linear, rates of convergence in expectation are derived that are of order $O(\log T/T)$ and $O(1/T^{1-\epsilon})$ for any $\epsilon > 0$, respectively under the availability of two-sample and one-sample oracles, respectively. Importantly, under the availability of the two-sample oracle, the aforementioned rate is actually agnostic to the relationship between the confounder and the instrumental variable, demonstrating the flexibility of the proposed approach in alleviating the need for explicit model assumptions required in recent works based on reformulating the problem as min-max optimization problems. Experimental validation is provided to demonstrate the advantages of the proposed algorithms over classical approaches like the 2SLS method.

CO286 Room Virtual R02 STATISTICAL MODELING AND INFERENCE UNDER COMPLEX STRUCTURE (VIRTUAL)
Chair: Jingru Mu
C0427: Estimation and inference of quantile spatially varying coefficient models over complicated domains

Presenter: **Myungjin Kim**, Kyungpook National University, Korea, South

Co-authors: Lily Wang, Huixia Judy Wang

The aim is to present a flexible quantile spatially varying coefficient model (QSVCM) for the regression analysis of spatial data. The proposed model enables researchers to assess the dependence of conditional quantiles of the response variable on covariates while accounting for spatial nonstationarity. The approach facilitates learning and interpreting heterogeneity in spatial data distributed over complex or irregular domains. A quantile regression method that uses bivariate penalized splines in triangulation is introduced to estimate unknown functional coefficients. The L2 convergence of the proposed estimators is established, demonstrating their optimal convergence rate under certain regularity conditions. An efficient optimization algorithm is developed using the alternating direction method of multipliers (ADMM). Wild residual bootstrap-based pointwise confidence intervals are developed for the QSVCM quantile coefficients. Furthermore, reliable conformal prediction intervals are constructed for the response variable using the proposed QSVCM. Simulation studies show the remarkable performance of the proposed methods. Lastly, the practical applicability of the methods is illustrated by analyzing the mortality dataset and the supplementary particulate matter (PM) dataset in the United States.

C0251: Exact statistical inference for some transformed gamma distributions: A generalized inference approach

Presenter: **Bowen Liu**, University of Missouri - Kansas City, United States

The transformed gamma distributions are frequently used for modeling extreme events across diverse fields, including hydrology, meteorology, and the insurance industry. Although these distributions demonstrate strong goodness-of-fit characteristics, exact inference for extreme quantiles remains challenging. Currently available parametric methods for quantile inference typically rely on approximation or bootstrapping techniques, which often demonstrate poor performance, particularly when sample sizes are relatively small. An exact inferential approach is presented for several transformed gamma distributions using generalized inference methodology. Additionally, the relationship is explored between the generalized inference method developed by Weerahandi and generalized fiducial inference. To evaluate the performance of the exact inference procedure,

a series of simulation studies are conducted. Results from these simulations indicate that the exact inference method consistently outperforms approximation-based and bootstrapping methods across multiple scenarios. Moreover, the practical applicability of the method was validated using multiple real-world datasets.

CO278: **Decoding spatial tissue architecture: A scalable Bayesian topic model for multiplexed imaging analysis**

Presenter: **Xiyu Peng**, Texas A&M University, United States

Recent progress in multiplexed tissue imaging is advancing the study of tumor microenvironments to enhance the understanding of treatment response and disease progression. Despite its popularity, there are significant challenges in data analysis, including high computational demands that limit feasibility for large-scale applications and the lack of a principled strategy for integrative analysis across images. To overcome these challenges, a spatial topic model is introduced, designed to decode high-level spatial architecture across multiplexed tissue images. The method integrates both cell type and spatial information within a topic modeling framework, originally developed for natural language processing and adapted for computer vision. Its performance is benchmarked through various case studies using different single-cell spatial transcriptomic and proteomic imaging platforms across different tissue types. It is shown that the method runs significantly faster on large-scale image datasets, along with high precision and interpretability. It consistently identifies biologically and clinically significant spatial topics, such as tertiary lymphoid structures.

C1008: **A semiparametric model approach to data integration under missing not at random**

Presenter: **Danhyang Lee**, Southern Methodist University, United States

Co-authors: Jae Kwang Kim

While probability sampling has long been the foundation of valid population inference, practical challenges such as rising costs and declining response rates have increased the use of non-probability samples. However, these alternative data sources are susceptible to significant selection bias. Data integration, which combines a non-probability sample with a reference probability sample, offers a potential solution. Most existing methods rely on the strong assumption of ignorable selection, known as missing at random. Although recent approaches have been developed for non-ignorable selection mechanisms, they typically depend on restrictive parametric models for the selection process. A unified data integration framework that accommodates both ignorable and non-ignorable selection mechanisms is proposed. An estimator is introduced for finite-population means developed within an empirical-likelihood framework. To mitigate the risk of model misspecification, the approach models the non-probability sample selection mechanism using a flexible semiparametric propensity-score specification. The resulting estimators combine calibration weighting for enhanced efficiency with the semiparametric propensity scores to reduce selection bias. Furthermore, closed-form variance expressions are derived, eliminating the need for computationally intensive replication techniques.

C1124: **Robust Bayesian elastic net with spike-and-slab priors**

Presenter: **Xi Lu**, University of Houston, United States

In high-dimensional regression problems, the demand for robust variable selection arises due to the commonly observed outliers and heavy-tailed distributions of the response variable, as well as model misspecifications when structured sparsity is ignored. The robust elastic net in both the frequentist and Bayesian frameworks has received much attention in recent years for the robust identification of important omics features. A robust Bayesian elastic net with spike-and-slab priors is proposed, which overcomes the major limitations of the existing family of elastic net methods. Specifically, a fully Bayesian method is developed that builds on the robust likelihood function to safeguard against the heterogeneity of complex diseases while accounting for structured sparsity. Incorporation of the spike-and-slab priors in the Bayesian hierarchical model has significantly improved accuracy in shrinkage estimation and variable selection. The advantages of the proposed method have been demonstrated through the simulation study of data with independent and identically distributed random errors and heterogeneous random errors over multiple versions of elastic net regularization methods and other alternatives. The analysis of SNP data with strong LDs from the Nurse Health Study (NHS) has also revealed the superiority of the proposed method.

CO080 Room BCB 206 NONLINEARITIES AND APPLICATIONS: A TRIBUTE TO H. PESARAN (VIRTUAL)

Chair: Willi Semmler

CO694: **Dynamic time-varying betas and climate and political risk factors for four largest Australian banks**

Presenter: **Timo Terasvirta**, Aarhus University, Denmark

Co-authors: Annastiina Silvennoinen, Glen Wade

Two types of risk are considered - the climate risk and the political one - for four major Australian banks, commonly called the Big Four. They account for about 20% of the value of the Australian stock market and are therefore important to the Australian economy. A new time-varying beta called the dynamic time-varying beta is first introduced. It is related to but different from (and less noisy than) the dynamic conditional beta that forms a part of a previous definition of the climate risk. Maximum likelihood estimators for this new beta are consistent and asymptotically normal. The dynamic time-varying beta is used for estimating the climate and the newly defined political risk for the Big Four using portfolios designed for Australian conditions. The results are compared with climate risk estimates based on the dynamic conditional beta and generated by V-Lab, a website of the New York University, Stern School of Business. While the climate risk estimates do not suggest a positive climate risk for the Big Four, the V-Lab estimates do. These differences are to a large extent due to different climate portfolios, probably less so to different betas.

CO740: **Co-extremal shocks and VAR analysis**

Presenter: **Stefan Mittnik**, Ludwig Maximilians University Munich, Germany

There is substantial empirical evidence that macroeconomic and financial shocks deviate from Gaussianity. A past study offered a theoretical explanation for this phenomenon by demonstrating that aggregation through production networks can transform micro disturbances into non-Gaussian macro shocks. Other studies cautioned that neglecting tails' properties may bias inference. In light of this, vector autoregressive (VAR) models have been refined to reflect this property by allowing errors to follow a multivariate Student's t-distribution, a Laplace distribution, or a generalized hyperbolic distribution. Based on the assumption of local ellipticity, the use of quantile-implied tail correlations is proposed as a natural approach to capturing non-Gaussianity and, in particular, to assessing the implications of extreme shocks in VAR modeling. When applied within a Global VAR (GVAR) framework, this approach can sharpen cross-country spillover analysis by better capturing the occurrence and implications of joint extremes in international shocks. This could lead to more realistic risk assessments and deeper, crisis-relevant policy insights at the global level.

CO848: **Technology diffusion and CO2 emission reduction**

Presenter: **Pu Chen**, Curtin University, Australia

Co-authors: Willi Semmler

The aim is to investigate the role of technology in the global green transition through both theoretical and empirical lenses. A dynamic general equilibrium (DGE) model is developed, based on the regional framework of a prior study, to examine how technological progress and international spillovers contribute to reducing CO emissions. In the model, technology influences both productivity and emissions intensity, with regional heterogeneity in innovation capacity and spillover absorption affecting the dynamics of transition. The analysis highlights the potential for cooperative technology strategies to accelerate decarbonization and enhance global welfare. Guided by the theoretical framework, the impact of technology is empirically assessed on emissions using a Global Vector Autoregression (GVAR) model. Employing panel data on CO emissions, technological indicators, and macroeconomic variables across a wide set of countries, the international transmission of technology shocks and their effects on emission trajectories are quantified. The empirical findings provide robust support for the hypothesis that cross-border technology diffusion

significantly contributes to emission reductions. These results underscore the importance of international cooperation and investment in clean technologies as central components of effective global climate policy

C1016: **Dynamic impact of decarbonization budget neutral fiscal policy on preferences, output and employment**

Presenter: **Feridoon Koohi-Kamali**, New School for Social Research, United States

Co-authors: Willi Semmler

Neutral climate taxation, financing an equal amount of fossil fuel reduction, is a relatively neglected tool of environmental public policy, which we address here. A model-guided approach is first presented to study budget-neutral climate policy impacts on preferences based on equal price elasticities between high and low carbon-intensive sectors. Second, a panel data method is empirically used to consider long-run homogeneity for an assessment of a balanced budget CO₂ diminishing policy on employment and output employing a demeaned tax subsidy formulation for a balanced budget variable that identifies periods of increasing and decreasing CO₂ emission by respectively negative and positive CO₂ log differences periods. Applied to a sample of 15 OECD economies over 1990-2023, the pooled mean group estimator has a better performance compared with other panel data models. Its long-run homogeneity assumption cannot be rejected, and policy coefficient estimates are significantly positive. The tests of the baseline model with cyclical time trend select the pooled mean group as the preferred model against the mean group and dynamic first differenced models. The model is stable, stationary, and robust to country-specific effects. The empirical outcome also confirms the model-guided approach that predicts the emission-reducing impact of budget-neutral fiscal policy on the energy transition.

CO021 Room BCB 207 INNOVATIONS IN FINANCE AND INSURANCE

Chair: Edit Rroji

C0349: **The role of market and mortality jumps in optimal DC pension schemes**

Presenter: **Immacolata Oliva**, Sapienza University of Rome, Italy

Co-authors: Davide Feleppa

The demographic shift marked by longer lifespans and lower birth rates poses challenges for worldwide PAYG pension systems. This shift has led to increased interest in pre-funded schemes, where individuals manage their investments throughout their careers to ensure sufficient retirement wealth. In addition, sudden events like pandemics or financial crises amplify uncertainty. The COVID-19 pandemic caused 1.6 million excess deaths in Europe (+8%), impacting pension and insurance systems. Similarly, financial crashes underline market vulnerability. To address such risks, jump-diffusion models are adopted for both asset returns and mortality, capturing discontinuous shocks and heavy-tailed distributions. Optimal asset allocation is studied in a DC pension scheme using a mean-CVaR approach to manage downside risk. The model features a risk-free asset and a risky one with jump-diffusion dynamics. Longevity follows a mean-reverting process with jumps. Independent Poisson processes drive jumps in asset and mortality dynamics. A mortality-linked derivative is introduced for hedging. The investors' problem is framed as a dynamic stochastic optimization, solved by minimizing CVaR and maximizing terminal wealth. Semi-Lagrangian schemes are used for the inner problem, and gradient descent for the outer, based on viscosity solutions to a PIDE. Time-inconsistent preferences and pre-commitment strategies are included, and robust tools are offered for retirement planning under extreme risk.

C0351: **Actuarial perspectives to the study of the life care reverse mortgage**

Presenter: **Giovanna Apicella**, University of Udine, Italy

Co-authors: Emilia Di Lorenzo, Giulia Magni, Marilena Sibillo

Reverse mortgage and long-term care insurance architectures are combined within the same insurance product, called Life Care Reverse Mortgage (henceforth LCRM). LCRM redistributes the borrower's payouts along the policy life, since it allows the borrower to shift into a new regime characterized by an enhanced benefit, after the inception of not-self-sufficiency. While not changing substantially the insurer's financial exposure, LCRM allows a guided planning of retirement savings. Actuarial methods are applied, and a formal mathematical structure is devised to build a time-dependent profit/loss function that describes over time the LCRM contract's value for the lender. In an example of numerical application, framed in the Italian context, ex-ante is monetized, the profit-loss function for a pool of LCRM contracts, showing the performance of this novel hybrid insurance contract against the typical reverse mortgage contract.

C0361: **Robust selection with carbon penalty**

Presenter: **Ilaria Stefani**, University of Parma, Italy

Co-authors: Immacolata Oliva

The aim is to present a robust, dynamic optimization model for continuous-time portfolio allocation, considering an investor with access to both green and brown assets in the stock market. It is assumed that the stock price dynamics are governed by a stochastic volatility model, where the instantaneous precision is modeled. The concept of 'penalized wealth' is introduced, where the wealth of an investor is adjusted for the losses incurred when investing in brown assets, which are assumed to bear higher risks, with the penalty being dependent on a parameter related to carbon emissions. By solving a system of ordinary differential equations (ODEs), a closed-form solution is derived to the optimization problem. Theoretical results highlight the significant influence of ambiguity aversion parameters on optimal policies. Furthermore, the parameter governing carbon emissions acts as an additional risk aversion, so the higher the emissions, the lower the portfolio allocation in the asset brown. A comprehensive numerical analysis, conducted with real data, validates these findings. Finally, the importance of ambiguity aversion is emphasized in investment decisions, particularly in the context of environmental sustainability and climate-related financial risks.

C0559: **A SVJJ model based on a compound CARMA(p,q)-Hawkes process**

Presenter: **Andrea Perchiazzo**, University of Milan, Italy

Co-authors: Lorenzo Mercuri, Edit Rroji

A stochastic volatility model is introduced with correlated jumps, incorporating a self-exciting effect in the intensity dynamics. As a primary goal, a pricing formula is derived based on the compound CARMA(p,q)-Hawkes framework, where the stochastic volatility is influenced by the quadratic variation of the counting process in the log-price dynamics. Additionally, a simulation algorithm is constructed for the jump term founded on the thinning algorithm. This algorithm is rooted in the existence of a Hawkes intensity with an exponential kernel, which serves as an upper bound for the CARMA(p,q)-Hawkes intensity. Finally, numerical and empirical analyses are presented.

CO250 Room BCB 208 PRICES AND WAGES INFLATION

Chair: Matteo Luciani

C0335: **Fiscal policy and inflation in the Euro area**

Presenter: **Lorenzo Mori**, University of Padova and Bank of Italy, Italy

Co-authors: Guido Ascari, Dennis Bonam, Andra Smadu

The relationship between fiscal policy and inflation dynamics in the Euro Area is investigated, with a focus on the post-pandemic inflation surge. Using a BVAR identified via sign restrictions, the effects of various demand and supply side shocks are disentangled, including fiscal policy, on inflation. First, while both positive demand and adverse supply shocks contributed to the inflation surge, demand shocks were relatively more important. Second, fiscal stimulus played a substantial and progressively increasing role, particularly in influencing domestic-based measures of inflation. Finally, the relative impact of fiscal shocks on inflation dynamics varies across (selected) Euro area countries.

C0513: **Lessons from the co-movement of inflation around the world**

Presenter: **Danilo Cascardi-Garcia**, Federal Reserve Board, United States

Co-authors: Luca Guerrieri, Matteo Iacoviello, Michele Modugno

The aftermath of the COVID-19 pandemic saw a synchronized surge in inflation around the world. However, the nature of this co-movement is ambiguous at least in three dimensions. First, whether the co-movement is driven by the trend component of inflation, or also by its cyclical part. Second, whether this co-movement is a phenomenon concerning sectors more exposed to the international markets' dynamics, or whether it percolates to other sectors. Third, what are the drivers of this co-movement. These ambiguities are resolved by looking at the behavior of core and noncore components of inflation across 20 different countries through the lenses of a rich model that divides each component into four elements: A common trend, a common cyclical component, an idiosyncratic trend, and an idiosyncratic cyclical component. It is found that the global component of inflation accounts for a large part of the variation in world inflation, including during the inflationary episodes of the energy crises of the 1970s-80s, as well as the post-COVID period. Moreover, the recent rise in global inflation is linked to both commodity-driven shocks and to global activity shocks such as the sharp decline and rebound in activity due to the COVID-19 pandemic and the associated labor market shortages.

C0542: The drivers of post-pandemic inflation

Presenter: **Giorgio Primiceri**, Northwestern University, United States

Co-authors: Domenico Giannone

Post-covid inflation was predominantly driven by unexpectedly strong demand forces, not only in the United States, but also in the Euro Area. These forces resulted from a combination of surprisingly robust pent-up demand following the pandemic restrictions, exceptionally expansionary fiscal policies, and an unusually accommodative monetary stance by the Federal Reserve and the ECB. This monetary accommodation has mitigated the recessionary effects of adverse supply shocks and supported the recovery, though it has come at the cost of higher inflation.

C0574: Inflation and the gender wage gap: The role of belief frictions

Presenter: **Nicolo Maffei Faccioli**, Norges Bank, Norway

How does the gender wage gap (GWG) respond to inflation? It is found that the GWG widens following both supply- and demand-driven inflationary shocks after controlling for individual characteristics and industry and occupation sorting. This widening is driven by gender differences in labor market perceptions: during inflationary periods, regardless of the source, women interpret inflation as a sign of strongly deteriorating labor market conditions (i.e., as supply shocks) and men as a sign of mildly improving labor market conditions (i.e., as demand shocks). These differing perceptions reduce women's reservation wages and lead to lower nominal wage growth relative to men. To capture this mechanism, a New Keynesian model is developed with two types of workers, men and women, and search and match frictions, where neither worker observes the true shocks. Instead, each forms potentially biased beliefs about the shock. The model successfully replicates the cyclicity of the GWG observed in the data.

CO169 Room BCB 210 DEPENDENCE MODELING IN ACTUARIAL SCIENCE

Chair: Etienne Marceau

C1133: Parametric estimation of conditional Archimedean copula generators for censored data

Presenter: **Marie Michaelides**, Heriot-Watt University, United Kingdom

Co-authors: Helene Cossette, Mathieu Pigeon

The aim is to propose a novel approach for estimating Archimedean copula generators in a conditional setting by incorporating endogenous variables directly into the generator function. Traditional copula models often rely on the simplifying assumption that the dependence structure remains fixed across all values of the covariates. The method relaxes this assumption by allowing both the strength and the shape of dependence to vary with covariates. In addition, to pinpoint the levels of a continuous risk factor where the dependence structure undergoes significant changes, an iterative splitting algorithm is introduced, which identifies optimal splitting points in the covariate space, that is, the range of possible values of a covariate. The effectiveness of the methodology is demonstrated through applications in two diverse settings, a diabetic retinopathy study and a claims reserving analysis, showing that accounting for covariate influences leads to a more accurate capture of the underlying dependence structure, thereby enhancing the applicability of copula models, for example, in medical or actuarial contexts.

C1135: Extremal negative dependence and the strongly Rayleigh property-part A

Presenter: **Patrizia Semeraro**, Politecnico di Torino, Italy

Co-authors: Alessandro Mutti, Helene Cossette, Etienne Marceau, Alessandro Mutti

Negative dependence properties of multidimensional Bernoulli distributions are important in many areas of probability. However, the theory of negative dependence is more challenging than that of positive dependence, and there are still open problems. While extremal positive dependence is an important concept clearly defined by the notion of comonotonicity, extremal negative dependence is still an open issue. Extremal negative dependence is characterized in Frechet classes of multidimensional Bernoulli distributions by completing the theory of negative dependence with the theory of dependence orders. Specifically, we prove that in the class of multidimensional Bernoulli distributions, extremal negative dependence is properly represented by the notion of sigma-countermonotonicity. Indeed, building on the geometrical representation of the class of Bernoulli variables, it is proven that the class of sigma-countermonotonic distributions is a convex polytope in the simplex of multidimensional Bernoulli distributions, and it is proven that it is an antichain that satisfies some minimality conditions with respect to the strongest negative dependence orders.

C1137: Extremal negative dependence and the strongly Rayleigh property-part B

Presenter: **Alessandro Mutti**, Politecnico di Torino, Italy

Co-authors: Helene Cossette, Etienne Marceau, Patrizia Semeraro

Multivariate Bernoulli distributions are essential in the modeling of binary data in a wide variety of contexts, such as actuarial science, quantitative risk management, machine learning, natural language processing, and bioinformatics. The aim is to present the second part of the investigation about extremal negative dependence for Bernoulli random vectors and the strongly Rayleigh property. The strongly Rayleigh property is the strongest negative dependence notion, implying negative association. The strongly Rayleigh property is often defined using the probability generating function of the Bernoulli random vector. Notably, a Bernoulli random vector satisfies the strongly Rayleigh property if its probability generating function is a stable polynomial. Methods to construct families of multivariate Bernoulli distributions that satisfy both the strongly Rayleigh property and the Sigma-countermonotonicity property (discussed in the first part of the investigation) are presented. Those notions are explained using the class of conditional Bernoulli distributions. Both the strongly Rayleigh property and the Sigma-countermonotonicity property are illustrated through examples in actuarial science and risk management.

C1179: A zero-inflated mixed-effects spatial point process for grouped storm loss data

Presenter: **Lisa Gao**, University of Waterloo, Canada

Co-authors: Sebastien Jessup

The increasing granularity of third-party weather and exposure information can allow insurers to more effectively predict weather-related losses. However, loss outcomes are often reported in spatially grouped observations, such as at the postal code level, so higher resolution predictors are aggregated to align with the granularity of the outcome in standard analyses. Assuming an underlying zero-inflated mixed-effects spatial point process framework for claims arising from a common storm, we derive a model for unbalanced, zero-inflated multivariate count data that incorporates rich weather and exposure predictors observed at different levels of spatial granularity to predict claim patterns. The model accommodates the dependence between locations affected by a common storm in the excess zeros, as well as in the joint claim counts. Using real property exposure

and loss data, the role of spatial dependence and granular predictors is highlighted in predicting localized storm losses.

CO306 Room BCB 212 ADVANCES IN ASSET PRICING
Chair: Sicong Li
C0496: Textual analysis of short-seller research reports

Presenter: **Xiao Han**, Bayes Business School, City St Georges, University of London, United Kingdom

Co-authors: Jules Van Binsbergen, Alejandro Lopez Lira

Using survey cash-flow expectations, it is found that investors underreact to the negative cash-flow news included in short-seller reports. On average, target firms earn abnormal returns of -4.9% on the publication day, and subsequent price revisions equal -15% over the next 12 months. A novel text-based fraud measure is introduced, and it is found that reports more related to fraud predict larger negative long-term abnormal returns. Using large language models, it is also found that other allegations, such as claims of overvaluation, do not have predictive power. Furthermore, short-seller research reports predict significant reductions in future real investment and stock issuances, with those heavily emphasizing fraud predicting larger declines. A model featuring expectation stickiness combined with limits to arbitrage is consistent with the long-term price impacts of short-seller reports.

C0768: Rethinking mutual fund performance: From traditional alpha to achievable alpha

Presenter: **Alberto Martin-Utrera**, Iowa State University, United States

Co-authors: Victor DeMiguel, Raman Uppal

Mutual-fund performance is traditionally evaluated using alpha, which measures the utility gain of an unconstrained investor who has access to the fund in addition to the benchmark factors. It is proven that the utility gain of shortsale-constrained investors is instead measured by achievable alpha, estimated using only those factors with strictly positive weight in the shortsale-constrained benchmark-factor portfolio. Empirically, active-fund management is less valuable for constrained investors: While 57.42% of funds have positive traditional gross alpha, only 30.77% have positive achievable gross alpha for a benchmark containing eight Vanguard funds. Achievable alphas predict fund flows, particularly during market turmoil.

C0771: Trading volume alpha

Presenter: **Chao Zhang**, HKUST(GZ), China

Rather than focusing on predicting asset return moments, the economic benefits are modeled for predicting individual stock trading volume. Volume forecasts are translated into a component of expected trading costs, and their value is analyzed through a portfolio framework. By recasting the volume prediction problem into a portfolio optimization problem that trades off tracking error versus net-of-cost performance, volume predictions are quantified into economic outcomes. Incorporating the economic loss function directly into a machine learning algorithm yields better out-of-sample performance than commonly used statistical loss functions. While volume is only one component of what drives trading costs, it is highly predictable, readily available, and its economic benefits are as large as those from stock return predictability.

C0992: Statistical arbitrage without arbitrage

Presenter: **Paolo Zaffaroni**, Imperial College London, United Kingdom

Co-authors: Valentina Raponi

Statistical arbitrage strategies are quantitative trading strategies based on alpha, i.e., the signal extracted from the residuals when fitting an asset pricing model to return data. A normative theory is developed that combines the insights of mean-variance portfolio choice with statistical arbitrage in the no-arbitrage setting of the APT, specifically studying statistical arbitrage without arbitrage opportunities. A novel two-fund separation result is established, that combines the inefficient statistical arbitrage portfolio and betaportfolio (i.e., the portfolio stemming from factor asset pricing models), showing how their combination can span the efficient frontier. In general, as the number of assets increases, the statistical arbitrage portfolio dominates the β portfolio both in terms of the magnitude of its weights and in terms of its SR. Building on insights from mean-variance portfolio choice, it is demonstrated how to construct a special statistical arbitrage portfolio that does not require estimating alpha, contrary to the prevailing view. When statistical arbitrage is combined with factor asset pricing modelling, alpha and the factors premia might not be jointly identified, jeopardizing the possibility of constructing a statistical arbitrage portfolio.

CO339 Room BCB 213 ASYMPTOTICS IN INFERENCE AND COMPUTATION FOR STOCHASTIC PROCESSES
Chair: Nakahiro Yoshida
C0263: Analytical approximations for diffusion processes with asymptotic expansion

Presenter: **Emanuele Guidotti**, University of Lugano, Switzerland

Co-authors: Nakahiro Yoshida

A comprehensive framework is presented for constructing analytical approximations of functionals of diffusion processes using asymptotic expansions in small noise regimes. Building on the Malliavin calculus, high-order expansions are derived for the probability density function, cumulative distribution function, characteristic function, moments, and quantile function of general diffusion processes. The method systematically develops these expansions for joint, marginal, and conditional distributions. A formal computational scheme is proposed to compute expansion formulas of arbitrary order. This approach allows accurate approximations in models where direct analytical solutions are unavailable, and it is broadly applicable to problems in mathematical finance, filtering, and inference for stochastic systems.

C0317: An estimator for the estimation error and hypothesis testing in non-identifiable models, including machine learning

Presenter: **Junichiro Yoshida**, University of Tokyo, Graduate School of Mathematical Sciences, Japan

Co-authors: Nakahiro Yoshida

To improve the affinity between machine learning and classical statistical methods, it is important to analyze the estimation error (defined as the difference between the expected loss of the estimator and that of the true parameter value) in non-identifiable models, including machine learning. However, in complex non-identifiable models, the estimation error is known to be difficult to calculate specifically. To solve this problem, an estimator is proposed for the estimation error and its confidence interval that can be applied to non-identifiable models, which enables the combination of classical statistical methods, such as hypothesis testing, with machine learning.

C0511: Neural Hawkes: Non-parametric estimation in high dimension and causality analysis in cryptocurrency markets

Presenter: **Ioane Muni Toke**, CentraleSupélec, France

Co-authors: Timothee Fabre

A novel approach to marked Hawkes kernel inference is proposed, which is named the moment-based neural Hawkes estimation method. Hawkes processes are fully characterized by their first and second order statistics through a Fredholm integral equation of the second kind. Using recent advances in solving partial differential equations with physics-informed neural networks, a numerical procedure is provided to solve this integral equation in high dimensions. Together with an adapted training pipeline, a generic set of hyperparameters is given that produces robust results across a wide range of kernel shapes. An extensive numerical validation is conducted on simulated data. Two applications of the method are finally proposed for the analysis of the microstructure of cryptocurrency markets. In a first application, the influence of volume is extracted on the arrival rate of BTC-USD trades and in a second application, the causality relationships and their directions are analyzed amongst a universe of 15 cryptocurrency pairs in a centralized exchange.

C0545: Parameter estimation for a linear parabolic SPDE in two space dimensions with a small noise using spatiotemporal data

Presenter: **Masayuki Uchida**, The University of Osaka, Japan

Co-authors: Yozo Tonaki, Yusuke Kaino

The purpose is to study the estimation of unknown parameters in a second-order linear parabolic stochastic partial differential equation (SPDE) in two spatial dimensions, driven by a Q-Wiener process with a small noise, using high-frequency spatio-temporal observations. Previous works focused on minimum contrast estimators (MCEs) for unknown coefficients in a second-order linear parabolic SPDE in one space dimension driven by a cylindrical Wiener process based on high-frequency spatiotemporal data. Another study further developed parametric adaptive estimators for a second-order linear parabolic SPDE in one space dimension with small noise. The methodologies of prior studies are extended to a linear parabolic SPDE in two dimensions with a small noise, and MCEs are proposed for the diffusive and advective parameters using temporal and spatial increments. Furthermore, an estimator is introduced for the reaction parameter based on an approximate coordinate process. Simulation studies are presented to evaluate the performance of the proposed estimators.

CO277 Room BCB M201 MULTIVARIATE EXTREMES

Chair: Ana Ferreira

C0882: Generalized logistic extreme value distributions

Presenter: **John Nolan**, American University, United States

Multivariate Frechet laws are a class of extreme value distributions that exhibit heavy tails and directional dependence controlled by an angular measure. Multivariate generalized logistic laws are a recently described subclass that are dense in a certain sense. It is shown that these laws are related to positive multivariate sum stable laws, which gives a way to simulate from these laws. The corresponding angular measure density is described, and expressions for the density of the distribution are given.

C0788: Bootstrap for the vector tail empirical process

Presenter: **Jun-ichiro Fukuchi**, Gakushuin University, Japan

A prior study derived the limiting process of the bootstrapped tail empirical process and its integral functional. Their results are extended to the vector tail empirical process (the vector TEP) based on the sample of a multivariate population. First, it is proved that the vector TEP converges weakly to a centered Gaussian process if the two-dimensional marginal distribution function admits an upper tail copula. Second, the limit of the bootstrapped vector TEP is derived. As a result, Efron's bootstrap is shown to be consistent for integral functionals under the no-bias assumption. As an application, a method of selecting the components is considered, whose high quantiles are the larger than a given value, which rely on the consistency of the bootstrap.

C1020: Dimension reduction within the geometric framework for multivariate extremes

Presenter: **Jeongjin Lee**, Lancaster University, United Kingdom

The framework of geometric extremes relies on the convergence of scaled sample clouds onto a limit set, characterized by a gauge function whose shape determines the extremal dependence structure. Recent statistical methodologies for estimating the limit set provide flexibility in capturing complex dependence structure. However, existing approaches are limited to relatively low dimensions. The focus is on a statistical model comprising a truncated gamma model for the radial component, conditional on the angular component, and an angular model defined on the simplex, leading to compositional data. In high dimensions, a common dimension reduction approach is to transform the angular components and apply principal component analysis. However, this approach can inflate variation in minor components and fail to capture linear patterns in the original compositional scale. As an alternative, compositional data is analyzed, directly on its original scale, using recently proposed methods that construct a nested sequence of lower-dimensional simplices, analogous to reverse principal component analysis. The performance of these methods is evaluated and compared through simulation studies and real data applications.

C0856: Sparse clustering for extremes

Presenter: **Nicolas Meyer**, CNRS, France

Identifying patterns of similar extremal behavior remains a central challenge in extreme value analysis. In a multivariate setting, the objective is to detect directions in which extreme events tend to concentrate, while in time series analysis, the goal shifts to identifying consecutive time intervals exhibiting extremal characteristics. Standard clustering techniques often prove inadequate due to the high dimensionality of the data relative to the scarcity of extreme observations. To address this limitation, sparse clustering methods offer a principled way to reduce dimensionality while preserving relevant extremal structure. The aim is to present a methodology based on the Euclidean projection onto the simplex to induce sparsity, leading to the concept of sparse regular variation. This framework provides a natural extension of multivariate regular variation, allowing for the identification of clusters of variables that exhibit similar tail dependence features. These clusters are often low-dimensional, which allows for a better interpretation. The proposed approach is illustrated using environmental and financial data, highlighting its interpretability and practical relevance.

C1234: Estimation of probabilities associated with failure sets

Presenter: **Ana Ferreira**, IST-ID, Portugal

Co-authors: Laurens de Haan

Frequently, one is faced with the occurrence of events that combine several variables taking extreme values. This happens in climatic events, with e.g., a combination of strong winds and high temperatures, and financial markets with portfolios combining several assets and/or stock indices. The estimation of the occurrence of such failure sets is theoretically demanding. The theory is revisited to make it more tractable for applications. The problem is reviewed from theoretical and practical viewpoints.

CO298 Room BCB M202 COMPUTATIONAL METHODS IN FINANCIAL ECONOMICS

Chair: Paul Schneider

C0808: Moment-driven predictive control of mean-field collective dynamics

Presenter: **Chiara Segala**, Università della Svizzera Italiana, Switzerland

The aim is to present a feedback control strategy for large-scale systems of interacting agents governed by nonlinear dynamics. Starting from an agent-based model with nonlocal interactions, the mean-field limit is considered to derive a tractable control formulation. The approach combines linearization techniques and Riccati equations to compute suboptimal feedback laws embedded in a model predictive control framework. Controls are updated using macroscopic quantities, specifically the first and second moments of the agent distribution. This circumvents the complexity of solving high-dimensional Hamilton-Jacobi-Bellman equations. A feedback law is constructed for the linearized mean-field model and iteratively refined based on the evolving moments. The resulting control strategy is efficient, scalable, and robust to partial observations. Theoretical performance estimates guide the selection of linearization points and update frequency. Numerical results show the effectiveness of the method in stabilizing self-organizing behaviors, with applications to opinion dynamics and collective motion.

C0728: Model-free bounds for option prices in incomplete markets

Presenter: **Yannick Dillschneider**, University of Amsterdam, Netherlands

A methodology is developed to quantify the uncertainty about option prices that persists in (statically) incomplete option markets due to finite strike and maturity grids as well as bid-ask price quoting. Uncertainty measures are provided for European option prices and their slopes (i.e., tail probabilities), each derived from model-free bounds that are consistent with a finite sample of observed bid and ask option prices and elementary no-arbitrage constraints in a multi-period setting. Novel explicit expressions are obtained for these bounds, which are simple and efficient to compute. Moreover, a decomposition of the bounds and uncertainty measures into contributions of strike, maturity, and price quote incompleteness

is suggested. In an empirical analysis of S&P 500 index options, the empirical uncertainty about option prices is quantified. Sizable uncertainty levels are documented for option prices and tail probabilities. Findings challenge the common belief that such "option-implied" (tail) information can be measured at sufficiently high precision, which has implications for many applications that use this information.

C0816: Nonparametric portfolio choices with firm characteristics

Presenter: **Rohan Sen**, USI Lugano, Switzerland

The aim is to present a novel approach to constructing optimal mean-variance portfolio rules, particularly when the number of assets is proportional to the number of samples. The approach models the portfolio weight of each asset as a nonparametric nonlinear function of firm-level characteristics and utilizes an equivalent unconstrained regularized formulation of the original mean-variance optimization problem. The framework generalizes to minimum-variance portfolios with short-sale constraints. The resulting portfolio rule is straightforward to implement, computationally efficient, and flexible.

C1122: Understanding the effect of input space dimension on kernel learning rates

Presenter: **Enrico Bignozzi**, USI, Switzerland

In recent years, the statistical-learning literature has made significant strides in non-parametric learning. These advances, however, have often run in parallel with the econometric tradition, which typically accounts for temporal dependence and explicitly links input dimensionality to the rate of risk convergence. In many econometric settings, asymptotic convergence rates depend directly on the dimensionality of the regressors, whereas most machine-learning work introduces alternative complexity measures and assumes i.i.d. data. The aim is to close that gap by deriving convergence rates for kernel ridge regression that make the dependence on the number of input variables explicit and accommodate temporally dependent data. Smoothness conditions are imposed on the regression target and relate the resulting rates to the eigenvalue decay of the kernel matrix, which is itself tied to the dimensionality of the inputs, yielding new asymptotic bounds that quantify the curse of dimensionality in non-i.i.d. scenarios.

CO197 Room BCB 307 DATA DEPTH IN STATISTICS AND MACHINE LEARNING

Chair: Hyemin Yeon

C1318: Anomaly detection for functional data

Presenter: **Hyemin Yeon**, Kent State University, United States

Co-authors: Xiongtao Dai, Sara Lopez Pintado

Anomaly detection is a fundamental task in Statistics and Machine Learning, especially as an initial step in exploratory data analysis. Techniques such as boxplot for univariate data, bagplot for multivariate data, and functional boxplot for functional data have been widely used to identify outliers across various domains. The focus is on functional data, introducing a new approach for detecting complex outliers using a novel concept of halfspace depth. This method is particularly effective in identifying functional shape outliers, which are known to be difficult to detect with traditional techniques. Numerical studies show that the proposed method outperforms existing approaches in recognizing various types of functional anomalies. Its practicality is demonstrated through applications to real-world datasets.

C1322: On a notion of graph centrality based on L_1 data depth

Presenter: **Seungwoo Kang**, Seoul National University, Korea, South

Co-authors: Hee-Seok Oh

A new measure to assess the centrality of vertices in an undirected and connected graph is proposed. The proposed measure, L_1 centrality, can adequately handle graphs with weights assigned to vertices and edges. The aim is to provide tools for graphical and multiscale analysis based on the L_1 centrality. Specifically, the suggested analysis tools include the target plot, L_1 centrality-based neighborhood, and local L_1 centrality. Most importantly, the focus is closely associated with the concept of data depth for multivariate data, which allows for a wide range of practical applications of the proposed measure. The tools are demonstrated with two interesting examples: the Marvel Cinematic Universe movie network and the bill co-sponsorship network of the 21st National Assembly of South Korea.

C1311: Differentially private multivariate medians

Presenter: **Kelly Ramsay**, York University, Canada

Co-authors: Shojaeeddin Chenouri, Dylan Spicker, Aukosh Jagannath

Data depth functions provide the standard framework for estimating a multivariate median. Currently, private versions are needed to allow for privacy-preserving, robust inference. Differentially private versions of various depth-based medians are explored, which are based on exponential mechanisms. In particular, the focus is on the theoretical properties, developing a general, tight finite sample deviations bound which can be applied to many depth-based medians at once, provided that the depths satisfy a simple condition. A method for relaxing this condition is also presented.

C1186: Robust fuzzy clustering for high-dimensional multivariate time series with outlier detection

Presenter: **Ziling Ma**, King Abdullah University of Science and Technology (KAUST), Saudi Arabia

Co-authors: Angel Lopez Oriona, Ying Sun, Hernando Ombao

Fuzzy clustering offers a natural approach to modeling partial memberships. However, most existing algorithms were developed for static, low-dimensional data. Such methods struggle with multivariate time series (MTS), where challenges include temporal dependence, unequal sequence lengths, moderate-to-high dimensionality, and frequent contamination by noise or artifacts. To overcome these challenges, RFCPCA is introduced, a robust fuzzy clustering framework that is, to the best of knowledge, the only approach specifically tailored to MTS that integrates (i) membership-informed subspace learning through common principal component analysis, (ii) the ability to handle unequal lengths and moderately high dimensions, and (iii) robustness against outliers via trimming, exponential down-weighting, and a noise cluster. This unique combination allows RFCPCA to capture latent temporal structure while remaining stable under contamination and informative about ambiguous or atypical series. Simulated and real Electroencephalogram data demonstrate that RFCPCA not only improves clustering accuracy compared to related methods, but also provides a more reliable characterization of uncertainty and atypical structure in MTS.

CO314 Room BCB 310 STATISTICAL CHANGE-POINT DETECTION

Chair: Hoseung Song

C0411: A conformal prediction framework for network systems

Presenter: **Rui Luo**, City University of Hong Kong, Hong Kong

Urban monitoring and sensor networks require reliable uncertainty quantification, yet methods with formal statistical guarantees are scarce. This is addressed by integrating graph neural networks (GNNs) with conformal prediction to create a framework for trustworthy, uncertainty-aware network analysis. This report details the advances in three core applications: conformal node classification for anomaly detection, conformal edge regression for traffic forecasting, and conformal route planning for risk-sensitive navigation. The methods are validated on 15 real-world datasets, demonstrating robust performance across diverse network tasks. Finally, future applications are discussed in critical urban systems, including cybersecurity and intelligent transportation.

C0588: Estimation and inference of change points in functional regression time series

Presenter: **Haotian Xu**, Auburn University, United States

The purpose is to study the estimation and inference of change points under a functional linear regression model with changes in the slope function. A novel functional regression binary segmentation (FRBS) algorithm is presented, which is computationally efficient and achieves consistency

in multiple change point detection. This algorithm utilizes the predictive power of piece-wise constant functional linear regression models in the reproducing kernel Hilbert space framework. A refinement step is further proposed that improves the localization rate of the initial estimator output by FRBS, and asymptotic distributions of the refined estimators for two different regimes are derived, determined by the magnitude of a change. To facilitate the construction of confidence intervals for underlying change points based on the limiting distribution, a consistent block-type long-run variance estimator is proposed. The theoretical investigation accommodates temporal dependence and heavy tails in both the functional covariates and the measurement errors. Empirical performance of the method is demonstrated through extensive simulation studies and applications to financial and economic datasets.

C0627: A general framework for online change point detection in nonparametric regression

Presenter: **Carlos Misael Madrid Padilla**, Washington University in St Louis, United States

The aim is to present a general framework for online change point detection in nonparametric regression models. In this setting, data arrive sequentially as covariate response pairs, and the underlying regression function is allowed to change at an unknown point in time. At each time step, the procedure fits two estimators, one to past data and one to recent data, by minimizing empirical squared loss over a chosen function class, and raises an alarm when the discrepancy between these estimators exceeds a time-dependent threshold. This approach accommodates a broad range of function classes, including kernel smoothers, spline-based models, trend filtering methods, and deep neural networks, allowing the procedure to adapt to various structural assumptions and complexities. General theoretical guarantees are established for the proposed method, ensuring control of the false alarm probability and sharp bounds on the detection delay.

C0680: Robust simultaneous detection of multiple breaks

Presenter: **Peiyun Jiang**, Tokyo Metropolitan University, Japan

Co-authors: Yoshimasa Uematsu

A novel multiple testing methodology is proposed for detecting structural changes in each variable of high-dimensional time series data. The approach enables break tests in the trend functions of individual time series without requiring prior knowledge of whether the noise components are stationary or integrated. The new multiple testing procedure controls the false discovery rate (FDR) at a predetermined level and accounts for complex dependencies in both the time-series and cross-sectional dimensions. The proposed method provides a robust framework for detecting structural breaks in high-dimensional time series settings.

CO156 Room BCB 312 ADVANCES IN OPTIMAL DESIGN AND SUBSAMPLING (VIRTUAL)

Chair: Nicholas Rios

C0565: Design and analysis for constrained order-of-addition experiments

Presenter: **Xueru Zhang**, University of Science and Technology Beijing, China

In an order-of-addition (OofA) experiment, the sequence of m different components can significantly impact the experiment's response. In many OofA experiments, the components are subject to constraints, where certain orders are impossible. For example, in survey design and job scheduling, the components are often arranged into groups, and these groups of components must be placed in a fixed order. If two components are in different groups, their pairwise order is determined by the fixed order of their groups. Design and analysis are needed for these pairwise-group constrained OofA experiments. A new model is proposed to accommodate pairwise-group constraints. A model for mixed-pairwise constrained OofA experiments is also introduced, which allows one pair of components within each group to have a pre-determined pairwise order. It is proven that the full design, which uses all feasible orders exactly once, is D - and G -optimal under the proposed models. Systematic construction methods are used to find optimal fractional designs for pairwise-group and mixed-pairwise constrained OofA experiments. The proposed methods efficiently assess the impact of question order in a survey dataset, where participants answered generalized intelligence questions in a randomly assigned order under mixed-pairwise constraints.

C0579: Optimal subsampling for linear models with heteroscedasticity

Presenter: **Jiayi Zheng**, George Mason University, United States

The rapid growth of data availability in modern scientific research and computation has led to increasingly large datasets across various fields. Studying linear relationships through regression analysis remains a fundamental task. However, despite advances in computer hardware, preserving and processing massive datasets continues to pose significant challenges. Subsampling methods, which focus on selecting the most informative subset of data based on optimal criteria, have been developed to address this issue, with the information-based optimal subdata selection (IBOSS) algorithm being a notable example. Additionally, heteroskedasticity in datasets presents a growing challenge. Methodologies are discussed for subsampling from big data, with a detailed explanation of two algorithms, weighted IBOSS and approximate nearest neighbor simulated annealing (ANNSA), and their performance in the presence of heteroskedasticity. A case study using real-world data is also included to illustrate their application.

C1256: REX-SUB: A scalable subsampling strategy for modeling large spatial datasets

Presenter: **Seiyon Lee**, George Mason University, United States

Co-authors: Nicholas Rios

Recent advances in data collection technologies have led to the emergence of massive spatial datasets, with measurements obtained at millions of spatial locations. Geostatistical models typically employ Gaussian processes (GPs) to capture spatial dependence, but standard GP fitting becomes prohibitive at such scales. A promising solution is optimal subsampling, where a subset of locations is selected that optimizes a criterion. The aim is to propose a randomized exchange algorithm for subsampling (REX-SUB), which efficiently selects small subsamples that minimize prediction errors in the fitted spatial GP models. To further improve computational efficiency, a scalable Vecchia approximation is embedded to GP's joint likelihood, which takes advantage of sparsity in the precision matrix to enable fast inference on the selected subsamples. Through a simulation study and an application to a remotely sensed precipitable water dataset, it is shown that REX-SUB yields lower mean squared prediction errors and interval scores compared to competing subsampling strategies.

CO299 Room BCB 403 MODERN APPROACHES TO DISTRIBUTIONAL AND HIGH-DIMENSIONAL DATA ANALYSIS

Chair: Wanli Qiao

C0621: Distributional outcome regression via quantile functions

Presenter: **Rahul Ghosal**, University of South Carolina, United States

Modern clinical and epidemiological studies widely employ wearables to record parallel streams of real-time data on human physiology and behavior. With recent advances in distributional data analysis, these high-frequency data are now often treated as distributional observations, resulting in novel regression settings. Motivated by these modelling setups, we develop a distributional outcome regression via quantile functions (DORQF) that expands existing literature with three key contributions: i) handling both scalar and distributional predictors, ii) ensuring jointly monotone regression structure without enforcing monotonicity on individual functional regression coefficients, iii) providing statistical inference via asymptotic projection-based joint confidence bands and a statistical test of global significance to quantify uncertainty of the estimated functional regression coefficients. The method is motivated by and applied to Actiheart component of Baltimore Longitudinal Study of Aging that collected one week of minute-level heart rate (HR) and physical activity (PA) data on 781 older adults to gain deeper understanding of age-related changes in daily life heart rate reserve, defined as a distribution of daily HR, while accounting for daily distribution of physical activity, age, gender, and body composition. Intriguingly, the results provide novel insights into the epidemiology of daily life heart rate reserve.

C1361: Doubly robust estimation of causal effects for random object outcomes with continuous treatments*Presenter:* **Bing Li**, The Pennsylvania State University, United States*Co-authors:* Satarupa Bhattacharjee, Xiao Wu, Lingzhou Xue

Causal inference is central to statistics and scientific discovery, enabling researchers to identify cause-and-effect relationships beyond associations. While traditionally studied within Euclidean spaces, contemporary applications increasingly involve complex, non-Euclidean data structures that reside in abstract metric spaces. The purpose is to introduce a novel framework for causal inference with continuous treatments applied to non-Euclidean data. To address the challenges posed by the lack of linear structures, Hilbert space embeddings of the metric spaces are leveraged to facilitate Frechet mean estimation and causal effect mapping. The framework can accommodate moderately high-dimensional vector-valued confounders and derive efficient influence functions for estimation to ensure both robustness and interpretability. Rigorous asymptotic properties of the cross-fitted estimators are established, and conformal inference techniques for counterfactual outcome prediction are employed.

C1009: Analysis on dynamics of deep neural network with high-dimensional structure*Presenter:* **Masaaki Imaizumi**, The University of Tokyo, Japan

Several topics related to the connection between statistics, machine learning, and dynamical systems are introduced. The first topic concerns the learning of the linearly non-separable structure by a neural network with simultaneous training. The result shows that a two-layer neural network can learn a linearly non-separable function even when both layers are updated simultaneously. To establish this result, the fine-grained tracking of neuron variability is characterized. The second topic discusses precise dynamical analysis of a deep neural network with a realistic architecture. Although the existing precise dynamical analysis depends on the infinite-width limit, the analysis allows finite-width and resolves some its limitations. The state evolution approach is applied to describe Gaussian fluctuations in the first-layer values and deterministic transform in the subsequent layers.

C0671: Embedding distributional data*Presenter:* **Wanli Qiao**, George Mason University, United States*Co-authors:* Ery Arias-Castro

The purpose is to adapt concepts, methodology, and theory originally developed in the areas of multidimensional scaling and dimensionality reduction for Euclidean data to be applicable to distributional data. The focus is on classical scaling and Isomap, prototypical methods that have played important roles in these areas, and showcase their use in the context of distributional data analysis. In the process, the crucial role that the ambient metric plays is highlighted.

CO323 Room BCB 405 ADVANCES IN BAYESIAN GRAPHICAL MODEL INFERENCE**Chair: Pariya Behrouzi****C0342: A Bayesian approach for inference on mixed graphical models***Presenter:* **Marina Vannucci**, Rice University, United States

Mixed data refers to a type of data in which variables can be of multiple types, such as continuous, discrete, or categorical. This data is routinely collected in various fields, including healthcare and social sciences. A common goal in the analysis of such data is to identify dependence relationships between variables, for an understanding of their associations. A Bayesian pairwise graphical model is proposed that estimates conditional independencies between any type of data. A flexible modeling construction is implemented, which includes zero-inflated count data and can also handle missing data. It is shown that the model maintains both global and local Markov properties. A spike-and-slab prior is employed for the estimation of the graph, and an MCMC algorithm is implemented for posterior inference based on conditional likelihoods. Performances on simulated data are assessed, and results are compared with existing methods. Finally, an analysis of real data from adolescents diagnosed with an eating disorder is presented.

C0482: Bayesian edge selection in mixed Markov random fields for ordinal and continuous variables*Presenter:* **Maarten Marsman**, University of Amsterdam, Netherlands*Co-authors:* Alexandros Beskos

Markov random field (MRF) graphical models are widely used in network psychometrics to model the network structure of psychometric variables. While existing methods typically handle either continuous or ordinal data, many psychological datasets include both. A novel MRF model is introduced for mixed ordinal and continuous variables. The model generalizes the ordinal MRF proposed by a recent study and is a special case of the conditional Gaussian graphical model. To estimate the network structure, a Bayesian variable selection approach is proposed using spike-and-slab priors on the edge weights. This prior structure aligns with existing methods for discrete MRFs and enables hierarchical modeling of the network structure. However, two challenges arise. First, the edge weight matrix must lie within the space of positive definite matrices, resulting in a joint prior with an intractable normalizing constant. This is addressed using a Metropolis algorithm that samples directly from the constrained space. Second, the MRF likelihood itself is computationally intractable. To overcome this, approximation strategies are explored, including pseudolikelihood methods. It is planned to implement the methodology in the bgms R package and integrate it into the JASP software platform for applied researchers.

C0580: Scalable Bayesian structure learning for Gaussian graphical models using marginal pseudo-likelihood*Presenter:* **Reza Mohammadi**, University of Amsterdam, Netherlands*Co-authors:* Lucas Vogels, Marit Schoonhoven, Ilker Birbil

Bayesian methods for learning Gaussian graphical models provide a powerful framework for addressing model uncertainty and incorporating prior knowledge. However, their applicability is often constrained by the high computational cost of exploring the joint space of graphs and precision matrices, particularly in large-scale settings. To overcome this limitation, an MCMC-based method is introduced that integrates the birth-death and reversible jump algorithms with the marginal pseudo-likelihood approach. By operating directly in the graph space, the method eliminates the need for intractable normalizing constants and precision matrix sampling, significantly reducing computational cost while preserving accuracy. These algorithms efficiently explore the full posterior graph space, enabling comprehensive model uncertainty quantification and seamless integration of prior structural knowledge. Notably, they produce reliable results in under an hour on standard hardware, even for datasets with over 1000 variables. The theoretical properties of the method are established, and its effectiveness is validated through an extensive simulation study and applications to gene expression data. Results show that the proposed algorithms outperform state-of-the-art Bayesian methods in both computational efficiency and accuracy. The algorithms are implemented in C++ and R and are accessible via the R package BDgraph.

C0581: Signpost testing to navigate the parameter space of the Gaussian graphical model with high-dimensional data*Presenter:* **Wessel van Wieringen**, Amsterdam University Medical Centers, Netherlands

A hypothesis test is presented to guide the search for the precision parameter of the Gaussian graphical model (GGM) from high-dimensional data. This parameter is unknown, but often information on it is available from a different setting. The unknown and different settings are assumed to be represented quantitatively by a parameter value. The direction between the values is a signpost in the parameter space. The hypothesis test evaluates the signpost for the true parameter. The test's null and alternative hypotheses state that no or some, respectively, relevant information for the parameter can be found in the direction of the signpost. The line segment is parameterized between the two settings by a one-dimensional parameter. The test statistic optimizes, within the class of regularized precision matrix estimators, the loss over this line segment. It measures the relevance of the signpost's direction for the problem at hand. The test statistic's null distribution is constructed or approximated asymptotically.

The signpost test's power is evaluated, and it is compared to the likelihood ratio test. The GGM's precision matrix of a pathway is learnt in a low-prevalent breast cancer subtype. In addition to in-house gene expression data, external data on a more prevalent subtype is available. The latter comes to mimic federated learning under the EU's GDPR, as a parameter estimate. The signpost test assesses this estimate's relevance for the in-house estimation problem.

CO116 Room BCB 406 RECENT ADVANCES ON NETWORK AND MATRIX DATA ANALYSIS
Chair: Jesus Arroyo
C0934: Vertex alignment and changepoint localization in network time series

Presenter: **Zachary Lubberts**, University of Virginia, United States

Co-authors: Mohammad Sharifi Kiasari, Avanti Athreya, Carey Priebe, Tianyi Chen, Sijing Yu, Youngser Park, Vince Lyzinski

Existing methodology for changepoint localization in an evolving time series of networks generally relies on accurately prescribed vertex correspondence between network realizations at different times. However, such vertex alignments are often misspecified or even unknown. To understand the impact of vertex misalignment on inference for dynamic networks, two illustrative models are constructed for network evolution, each with a similar changepoint. Different techniques are compared for changepoint localization, ranging from the simple network statistic of average degree to the more involved and recently developed procedure of Euclidean mirrors. In one model, vertex misalignment causes comparatively little error, and in the other, it seriously impairs localization, although the Euclidean mirror procedure can nevertheless extract a meaningful signal. It is shown how misalignment between network realizations at different times can effectively weaken their underlying correlation, impeding inference procedures that rely on accurate inference of such correlation. Graph matching and optimal transport is discussed, both of which are potential mechanisms for mitigating errors from misalignment, but which may also fail to improve inference under certain models. Simulations are presented that illustrate these varying effects on approaches to localization.

C0853: Computationally optimal estimation and inference of common subspaces

Presenter: **Joshua Agterberg**, University of Illinois Urbana-Champaign, United States

The statistical and computational limits for the common subspace model are investigated, a model wherein one observes a collection of symmetric low-rank matrices perturbed by noise, where each low-rank matrix shares the same common subspace. First, an estimator is proposed based on Grassmannian gradient descent initialized via a spectral sum of squared matrices, and it is shown that it achieves the optimal sin Theta error under a strong signal-to-noise ratio (SNR) condition, and further evidence is given that this SNR condition is necessary for a polynomial-time estimator to exist. Next, estimation and inference are used for the sin Theta distance itself, and it is shown that the estimator achieves an asymptotically Gaussian distribution with a bias term that vanishes under a strong signal requirement. Based on this limiting result, confidence intervals are proposed, and it is shown that they are minimax optimal, though the resulting confidence intervals require knowledge of the SNR. Then, it is turned to designing adaptive confidence intervals for the sin Theta error, and that adaptivity is information-theoretically impossible unless the SNR is sufficiently strong. Consequently, the results unveil a novel phenomenon: despite the SNR being "above" the computational limit for estimation, adaptive statistical inference may still be information-theoretically impossible.

C1357: Spectral embeddings of correlation networks

Presenter: **Keith Levin**, University of Wisconsin, United States

In many applications, weighted networks are constructed based on time series data. Most typically, a time series is associated with each vertex, and edge weights are given by the correlations between time series. The result is a network with a dependency structure among the edges that violates the assumptions of most common network models. Nonetheless, it is common to apply network embedding methods to networks built from correlations. The aim is to show that this violation of assumptions is not critical. A setup is considered in which a collection of time series signals is observed subject to noise, and a network is constructed based on correlations between the noisy series. It is proven that, under suitable conditions, applying the adjacency spectral embedding to the network of correlations recovers a population embedding of the true time series. Additionally, it is shown that the resulting embedding encodes (up to orthogonal rotation) the Fourier coefficients of the true time series. This appears to be folklore among the signal processing community in the context of principle component analysis, but it is, to the best of knowledge, new to the networks literature.

C1352: Network signflip parallel analysis for selecting the embedding dimension

Presenter: **Joshua Cape**, University of Wisconsin, Madison, United States

Co-authors: David Hong

The purpose is to investigate the problem of selecting the embedding dimension for large heterogeneous networks that have weakly distinguishable community structure. For a broad family of embeddings based on normalized adjacency matrices, a novel spectral method is introduced, that compares the eigenvalues of the normalized adjacency matrix to those obtained after randomly signflipping its entries. The proposed method, called network signflip parallel analysis (NetFlipPA), is interpretable, simple to implement, data-driven, and does not require users to carefully tune parameters. For large random graphs arising from degree-corrected stochastic block models with weakly distinguishable community structure (and consequently, non-diverging eigenvalues), NetFlipPA provably recovers the spectral noise floor (i.e., the upper-edge of the eigenvalues corresponding to the noise component of the normalized adjacency matrix). NetFlipPA thus provides a statistically rigorous randomization-based method for selecting the embedding dimension by keeping the eigenvalues that rise above the recovered spectral noise floor. Compared to traditional cutoff-based methods, the data-driven threshold used in NetFlipPA is provably effective under milder assumptions on the node degree heterogeneity and the number of node communities. Main results rely on careful non-asymptotic perturbation analysis and leverage recent progress on local laws for nonhomogeneous Wigner-type random matrices.

CO076 Room BCB 408 BAYESIAN SEMI- AND NON-PARAMETRIC METHODS II
Chair: Beatrice Franzolini
C0609: Spatiotemporal clustering with Neyman-Scott processes

Presenter: **Yixin Wang**, University of Michigan, United States

Neyman-Scott processes (NSPs) are point process models that generate clusters of points in time or space. They are natural models for a wide range of phenomena, ranging from neural spike trains to document streams. The clustering property is achieved via a doubly stochastic formulation: first, a set of latent events is drawn from a Poisson process; then, each latent event generates a set of observed data points according to another Poisson process. This construction is similar to Bayesian nonparametric mixture models like the Dirichlet process mixture model (DPMM) in that the number of latent events (i.e., clusters) is a random variable, but the point process formulation makes the NSP especially well suited to modeling spatiotemporal data. While many specialized algorithms have been developed for DPMMs, comparatively fewer works have focused on inference in NSPs. Novel connections are presented between NSPs and DPMMs, with the key link being a third class of Bayesian mixture models called mixture of finite mixture models (MFMMs). Leveraging this connection, the standard collapsed Gibbs sampling algorithm is adapted for DPMMs to enable scalable Bayesian inference on NSP models. The potential of Neyman-Scott processes is demonstrated on a variety of applications, including sequence detection in neural spike trains and event detection in document streams.

C0813: On L^2 posterior contraction rates in Bayesian nonparametric regression models

Presenter: **Paul Rosa**, University of Cambridge, United Kingdom

A central result in Bayesian nonparametrics is that under assumptions on the prior, the data-generating distribution (assuming a true frequentist model) and a semi-metric $d(\cdot, \cdot)$ on the space of regression functions that satisfy the so called testing condition, the posterior contracts around the

true distribution with respect to $d(\cdot, \cdot)$, and the rate of contraction can be estimated. In the nonparametric regression setting, the semi-metric $d(\cdot, \cdot)$ is often taken to be the Hellinger distance or the empirical L^2 norm (i.e., the L^2 norm with respect to the empirical distribution of the design) in the present regression context. However, extending contraction rates to the "integrated" L^2 norm usually requires more work and has previously been done, for instance, under sufficient smoothness or boundedness assumptions, which may not necessarily hold. It is shown that, for priors based on truncated random basis expansions and in the random design setting, a high probability two sided inequality between the empirical L^2 norm and the integrated L^2 norm holds in appropriate spaces of functions of low frequencies, under mild assumptions on the underlying basis (which can be for instance a Fourier, wavelet or Laplace eigenfunction basis), allowing us to directly deduce an L^2 contraction rate from an empirical L^2 one without further assumption on the true regression function.

C0279: A Stein gradient descent approach for doubly intractable distributions

Presenter: **Jaewoo Park**, Yonsei University, Korea, South

Bayesian inference for doubly intractable distributions is challenging because they include intractable terms, which are functions of parameters of interest. Although several alternatives have been developed for such models, they are computationally intensive due to repeated auxiliary variable simulations. A novel Monte Carlo Stein variational gradient descent (MC-SVGD) approach is proposed for inference for doubly intractable distributions. Through an efficient gradient approximation, the MC-SVGD approach rapidly transforms an arbitrary reference distribution to approximate the posterior distribution of interest, without necessitating any predefined variational distribution class for the posterior. Such a transport map is obtained by minimizing Kullback-Leibler divergence between the transformed and posterior distributions in a reproducing kernel Hilbert space (RKHS). The convergence rate of the proposed method is also investigated. The application of the method is illustrated to challenging examples, including a Potts model, an exponential random graph model, and a Conway-Maxwell-Poisson regression model. The proposed method achieves substantial computational gains over existing algorithms, while providing comparable inferential performance for the posterior distributions.

C0972: Model-based clustering of categorical data based on the Hamming distance

Presenter: **Lucia Paci**, Università Cattolica del Sacro Cuore, Italy

Co-authors: Raffaele Argiento

A model-based approach is developed for clustering categorical data with no natural ordering. The proposed method exploits the Hamming distance to define a family of probability mass functions to model the data. The elements of this family are then considered as kernels of a finite mixture model with an unknown number of components. Conjugate Bayesian inference has been derived for the parameters of the Hamming distribution model. The mixture is framed in a Bayesian nonparametric setting, and a trans-dimensional blocked Gibbs sampler is developed to provide full Bayesian inference on the number of clusters, their structure, and the group-specific parameters, facilitating the computation with respect to customary reversible jump algorithms. Extensions to overcome the independence assumption of the variables within the clusters are discussed. Model performances are assessed via a simulation study and reference datasets, showing improvements in clustering recovery over existing approaches.

CO108 Room BCB 409 STATISTICS AND SPORT

Chair: Brigitte Gelein

C0265: College football volatility: A Bayesian state-space model of the transfer portal and NIL impact

Presenter: **Ronald Yurko**, Carnegie Mellon University, United States

The landscape of American college football has changed dramatically in recent years with conference realignment, increased usage of the transfer portal, and the seismic legal ruling regarding athletes' name, image, and likeness (NIL). A Bayesian state-space model is used to capture the impact of transfers and NIL on the volatility of team strengths over time. Specifically, the classic autoregressive process for team strength is extended by modeling the between-season innovation variance as a function of incoming transfers and NIL rule changes. This approach enables predicting greater variance for team strengths in an upcoming season based on roster changes. Results from the playoff era (2014-2024) reveal variation between schools that have embraced the volatility of transfers (e.g., Deion Sanders at Colorado) and those who have been reluctant to change (e.g., Iowa and Clemson). Variability is explored in the effects between different positions of incoming transfers, and the predictive performance of the novel approach is compared with simpler state-space models.

C0384: Clustering of soccer possessions using a mixture of absorbing Markov point processes

Presenter: **Koffi Amezougui**, ENSAI (Ecole Nationale de la statistique et de l'analyse de l'information), France

Co-authors: Brigitte Gelein, Matthieu Marbac, Anthony Sorel

Motivated by the analysis and clustering of soccer game situations, a specific finite mixture of absorbing Markov point processes models is introduced to cluster soccer ball possessions based on their temporal, spatial, and categorical characteristics. Each possession is modeled as a marked point process that ends with a specific event that causes this loss of ball possession, such as a foul or an interception by an opposing player. Event types are modeled using finite-state Markov chains, inter-event times are assumed to follow Gamma distributions, and spatial displacements are represented by time-scaled truncated Gaussian distributions. Model parameters are estimated using a generalized expectation-maximization (GEM) algorithm. The model is validated through numerical experiments and applied to event data from professional soccer matches. The results reveal interpretable clusters that reflect distinct tactical behaviors. Meaningful possession patterns are extracted, providing insights into game dynamics for performance analysis and training.

C0431: An explainability approach for identifying winning strategies in rugby union

Presenter: **Sebastien Dejean**, University of Toulouse, France

Co-authors: Arnaud Odet, Cristian Pasquaretta

Interest in predicting sports match outcomes has grown significantly. However, the utilization of predictive models in enhancing tactical team performance remains relatively limited. A methodology is proposed that combines machine learning and algorithm explainability techniques. The case study on rugby union unfolds in two phases. First, the best modeling approach is identified for the data by establishing a prediction model based on performance indicators observed during games. Then, SHapley Additive exPlanations (SHAP) values are used to interpret the predictions of this model. Findings serve three purposes: (i) from a global standpoint, identifying performance indicators that potentially determine match outcomes; (ii) from an aggregated point of view, highlighting strengths and weaknesses of any given team; and (iii) from a local perspective, offering technical staff diagnostic analyses of past games.

C0987: Modelling time-varying training loads to assess their impact on sports injuries

Presenter: **Dae Jin Lee**, IE University, Spain

A flexible modeling framework is proposed to analyze time-varying exposures and recurrent events in team sports injuries. This framework leverages a piece-wise exponential additive mixed model to capture the cumulative and potentially complex effects of past exposures, such as high-intensity training loads, on current injury risk. To determine the optimal time window during which past exposures influence risk, a penalty-based approach is introduced. A simulation study is conducted to assess the performance of the proposed model under various scenarios, including different underlying weight functions and varying levels of heterogeneity across recurrent events. Finally, the application of this approach is demonstrated through a case study involving an elite male football team competing in Spain's LaLiga. The cohort includes time-loss injury data and external training load measures tracked via global positioning system (GPS) devices over the 2017-2018 and 2018-2019 seasons.

CC391 Room BCB 209 TEXT DATA IN ECONOMICS AND FINANCE**Chair: Bettina Gruen****C0307: How source, topic, & recency influence LLM-based sentiment predictions: Evidence from cross-sectional return prediction****Presenter:** Jule Schuettler, University of St.Gallen, Switzerland

The purpose is to investigate the predictive value of sentiment extracted from diverse financial text sources using state-of-the-art large language models (LLMs). In addition, it is examined whether incorporating topic information can enhance sentiment classification. The analysis is based on a comprehensive dataset comprising earnings call transcripts, newspaper headlines, and social media tweets. Using fine-tuned DeBERTa models, long/short portfolios are constructed based on sentiment predictions, and their trading performance is evaluated across sources. It is found that, while all sources contain predictive information, only tweets and headlines yield profitable strategies once trading costs are considered. Furthermore, it is demonstrated that combining multiple sources improves predictive performance by reducing portfolio volatility. A novel topic-aware LLM that incorporates topic information through unsupervised topic modeling is also introduced, and trading strategies based on this model are shown to outperform those that do not account for topical context. Finally, it is found that a rolling window training approach significantly enhances model performance, underscoring the importance of recency in dynamic financial environments.

C1393: Language shift in Polish and German economic journals: A corpus-based comparative analysis**Presenter:** Viktoriia Naboka-Krell, Justus-Liebig-University Giessen, Germany**Co-authors:** Anna Staszewska-Bystrova, Victor Bystrov, Peter Winker

English is widely recognized as the dominant language of scientific publications. The shift from national languages to English presents new opportunities and challenges for researchers, scientific institutions, and publishers. The results of a comparative analysis of the language shift are reported in two economic journals published in Poland and Germany. Specifically, a long period from 2000 to 2023 is considered. Using structural topic modeling, changes in the structures and dynamics of topics discussed in the journals are compared conditionally on the language composition of publications. Both topical structure/composition are looked at over time, as well as terminology usage within the topics. Using topic matching, topics common to both corpora are identified. Topics that were specific to one of the countries considered are also investigated, and possible reasons for the appearance or disappearance of these topics are discussed, i.e., changes in the relative importance over time.

C1469: ChatGPT-extracted news sentiment and the cross-section of U.S. stock returns**Presenter:** Lukas Petrasek, Charles University Prague, Czech Republic

The purpose is to analyze more than one hundred thousand news headlines to investigate their impact on U.S. equity prices, leveraging ChatGPT to extract sentiment, novelty, importance, and other key linguistic metrics. For every day between 1st January 2004 and 30th November 2024, an index ranging between -1 (bad) and 1 (good) is constructed, describing the overall sentiment of news published on that given day. Individual U.S. stock returns are then regressed on the sentiment index to determine their sensitivity to the sentiment of the news headlines. Based on these exposures, a portfolio-sorting approach is implemented to analyze the cross-sectional effects of news-driven sentiment on stock performance. Findings reveal that stocks that have been highly exposed to news sentiment in the previous year earn 3.7% higher annualized returns in the following month compared to stocks that have been less exposed to the sentiment index. This research contributes to AI-driven asset pricing models and offers practical implications for investors incorporating news-based factors into trading strategies.

C1228: Topic-based expected sentiment values for market volatility forecasting**Presenter:** Agnesa Hovhannisyan, University of Salerno, Italy**Co-authors:** Alessandra Amendola, Francesco Audrino

Various quantitative and qualitative factors are used in the literature in attempt to enhance financial volatility forecasting. With the emergence of textual analysis tools, the broader source of qualitative information is examined for incorporation into the financial models to introduce additional information for volatility predictions. The aim is to build a temporal sentiment index defined as Expected Sentiment Value, which is derived from the estimated topic-document probability distribution from topic modelling, combined with the machine learning-based sentiment scores. DeepSeek-V3 is used to extract sentiment scores from the data used. News headlines and descriptions are used to examine whether this approach enhances the forecasting power of the Heterogenous Autoregressive (HAR) model. Using data from three news outlets (CNBC, Guardian, and Reuters) for the S&P500 Index between March 20, 2018 and July 20, 2020, minor improvements are shown over the benchmark HAR model, however no significant support is found for enhanced forecasting power.

CC379 Room BCB 308 MACHINE LEARNING**Chair: Sanjoy Sinha****C1214: KLAN: Kolmogorov-Lorentz-Arnold neural networks****Presenter:** Andrej Svetlosak, The University of Edinburgh, United Kingdom**Co-authors:** Miguel de Carvalho, Clemente Ferrer, Fengxiang He, Johnny Lee, William Wu

Kolmogorov-Arnold neural networks (KANs) have been investigated as an alternative to multilayer perceptron networks. Based on the Kolmogorov-Arnold (KA) representation theorem, weights are replaced with non-linear learnable activation functions. This gives KANs flexibility and interpretability. A novel network architecture based on Lorentz's version of the KA theorem is proposed, which is called Kolmogorov-Laurenz-Arnold networks (KLAN). Under certain conditions, the outer function of the KAN architecture can be replaced with a single function. The result is an equally powerful, but sparser network, requiring fewer parameters and less computational effort. The properties of the proposed methods are demonstrated in a series of real and simulated experiments, and their performance is compared to state-of-the-art versions of KAN. It is concluded by discussing the advantages and shortcomings of the proposed approach.

C1243: A statistical framework for characterizing the response distribution of a large language model**Presenter:** William Wu, The University of Edinburgh, United Kingdom**Co-authors:** Miguel de Carvalho

Large language models (LLMs) often produce a spectrum of plausible responses to a single prompt; little attention has been paid to characterizing the conditional distribution of their responses, given a fixed query. The aim is to present a framework for analyzing LLM empirical distributions by leveraging sentence embeddings and similarity-based metrics. Total Recall and the RowBERT score are introduced as interpretable measures of variability of the ensemble responses and global similarity per reply. Building on these, a probabilistic scoring framework that quantifies the likelihood of a candidate response generated from a given ensemble is developed. Through simulations and experiments on multiple LLMs and open-source datasets, it is demonstrated that the empirical distribution formed by RowBERT score is stable, interpretable, and highly responsive to key generation parameters like temperature. The proposed heuristics-guided methodology offers a practical and interpretable framework for characterizing empirical output distributions from LLMs, providing critical insights into model robustness, trustworthiness, and the inherent variability of generated output.

C1227: Assessing the fairness of stability for binary classifiers**Presenter:** Hiroe Seto, Kyoto Women's University, Japan**Co-authors:** Michio Yamamoto, Shin-ichi Mayekawa

In recent years, binary classifiers have increasingly been used to make important decisions, such as hiring or lending. However, there are concerns that the predictions of binary classifiers may lead to unfairness. Unfairness is defined as a worse impact on protected classes, such as women or ethnic minorities, than on non-protected classes. Previous research has developed methods to assess the fairness of prediction accuracy between

protected and non-protected classes. However, there has been no method to assess the fairness of stability, which refers to the degree of agreement between predictions made by a model trained on a different dataset. Using an algorithm with unfair stability results in greater variance in predictions for protected classes than for non-protected classes, posing a risk of exacerbating unfairness in society. DOSACA is therefore proposed as a metric that assesses unfairness based on the difference in stability between protected and non-protected classes. Furthermore, to verify the practical usefulness of the proposed evaluation method, real-world data analysis is conducted using several open data sets. The results of this research provide new evaluation criteria for the development of fair binary classifiers and have important implications, particularly for the ethical and legal aspects of predictive algorithms in the real world.

C1325: **Spatial information integration for analyzing cellular drug responses**

Presenter: **Mahiro Yamamoto**, Doshisha University, Japan

Co-authors: Masaaki Okabe, Hiroshi Yadohisa

Single-cell RNA sequencing (scRNA-seq) provides high-resolution gene expression profiles but lacks spatial context, limiting tissue-level analysis. Spatial transcriptomics (ST) preserves spatial information but cannot achieve single-cell resolution. While optimal transport methods can integrate ST spatial information with scRNA-seq data, drug responses analysis (perturbation analysis) remains challenging, where post-perturbation ST data is not readily obtainable. The aim is to propose a novel approach using neural optimal transport to propagate spatial information from pre-perturbation spatially-annotated scRNA-seq data to post-perturbation scRNA-seq data. The method models cellular state transitions by learning transport mappings between gene expression distributions via neural networks, enabling correspondence between unpaired pre- and post-perturbation cells. This framework addresses the critical gap in spatially-aware perturbation analysis by leveraging pre-perturbation ST-scRNA-seq integration as a foundation for spatial context propagation. The method enables comprehensive analysis of spatially-resolved gene expression patterns in response to drug responses. The contribution enables spatial context-aware analysis of drug responses where direct spatial measurement is technically challenging, opening new possibilities for understanding spatially dependent cellular responses to various experimental drug responses.

CC443 Room BCB 309 SHORT TALKS: CFE

Chair: Peter Pedroni

C0194: **Convergence in overqualification rates across EU Countries**

Presenter: **Theologos Pantelidis**, University of Macedonia, Greece

Co-authors: Maria Karantali

Overeducation or overqualification is an undesirable employment situation in which an individual's qualifications exceed job requirements, and the level of education required to perform the work is lower than his/her educational background. Overeducation spurred a wide array of policy debates and studies on skills mismatch that focus on the surplus of human capital. Overqualification rates are used to investigate beta and sigma convergence across EU countries over the period 2005 to 2023. Various econometric techniques, including club convergence ones, are employed to examine convergence for the whole period under scrutiny and for two sub-periods as well. The analysis is further expanded to explore differences in overqualification rates by gender and by age groups (namely, 20-34 and 35-64). The results provide mixed evidence and multiple convergence clubs, suggesting that further policies are required to manage overeducation within the EU countries.

C1506: **A nonparametric test for social spillovers and interference**

Presenter: **Abhimanyu Gupta**, Queen's University, Canada

Co-authors: Margherita Comola, Camila Comunello

A novel nonparametric test is proposed for cross-unit dependence occurring through peers' attributes (interference) or peers' outcomes (social spillovers). While the test is nonparametric, it is reminiscent of classical Lagrange multiplier-type diagnostic tests. Indeed, the test statistic requires only estimates under the null to be computed and therefore avoids nonparametric estimation altogether, and is shown to have a convenient asymptotic standard normal distribution. Empirical illustrations are presented to show that the test is effective at detecting cross-unit dependence arising in a nonparametric fashion that existing approaches might not detect effectively.

C1513: **Estimating points of structural correlation instabilities with latent factor models**

Presenter: **Sojin Lee**, CUNEF Universidad, Spain

Cross-correlation structures contain valuable information about the underlying linkages among variables and the channels of spillovers across cross-sectional units. Instabilities in these structures often signal structural changes within the system. We propose a novel method for detecting instabilities in cross-correlation structures using a latent factor model framework. We introduce a suitable object, the column space of the loading matrix (the factor space), to capture structural correlation changes while free from the inherent identification issue of the latent model. The resulting detection criterion is based on an intuitive distance measure between two factor spaces, integrating both the detection and localization of breakpoints. The factor space-based framework can also accommodate extension to online early-warning methods. Within this framework, we discuss how classical approaches in factor model analysis coherently connect with the probability of structural change, thereby offering a baseline criterion for sequential early-warning methods. In applications, our methods effectively detect instability points consistent with the development of the subprime mortgage crisis, as well as major policy changes such as the repeal of the Glass-Steagall Act and the U.S.-China trade war.

C1518: **Measuring excess volatility & asset mispricings through the chiarella model**

Presenter: **Jutta Kurth**, Ecole Polytechnique, France

Co-authors: Jean-Philippe Bouchaud, Adam Aleksander Majewski

The two most well-accepted anomalies pervading all financial markets are trend/momentum and value. Their coexistence and interaction may be studied analytically and numerically via a stochastic dynamical system known as Chiarella's model and its extensions. This framework is generalized to include arbitrary drifts in the fundamental value and a trend signal that depends on the mispricing of the asset rather than its bare log-returns, while keeping the same linear stability. The fundamental value is a direct output of the calibration and does not rely on economic pricing models. Numerically, results on synthetic data as well as on several contracts from five different asset classes show that the model provides better fits. Instead of calibrating on asset classes, it is possible to calibrate on individual price series, showcasing a powerful improvement in the fitting procedure. Stylized facts are confirmed, e.g., bimodalities in the mispricing distribution (from a phenomenological P-bifurcation in the model), suggesting that prices are more often mispriced than correctly priced, as well as the non-monotonic relation between past trends and future returns for all asset classes. The robustness of results against perturbations in each parameter and their hierarchy is studied via sloppiness analysis. The excess volatility puzzle is quantified for all five asset classes by comparing the calibrated noise strength of the noise trader with that of the fundamental value.

C1522: **Likelihood specification testing in nonlinear panel data models with fixed effects**

Presenter: **Martin Schumann**, Maastricht University, Netherlands

Consistent maximum likelihood estimation typically depends on the correct specification of the likelihood. If the presence of individual heterogeneity is not accounted for, estimation and inference will frequently yield misleading results. We develop an information-matrix test for likelihood-based specification testing in nonlinear panel-data models with unobserved fixed effects. We show that allowing for fixed effects leads to an incidental parameters problem, rendering the information matrix test inconsistent if the number of time periods does not grow faster than the number of individuals. The test's incidental parameters bias is shown to depend on the score and the information bias of the pseudo-likelihood

on which the test statistic is based. We derive a bias-corrected version of the information matrix test and analyze its small-sample properties in a simulation study.

C1528: **Comparison between the Markowitz and Konno-Yamazaki Models for Portfolio Selection**

Presenter: **Jose Luis Esteves dos Santos**, Top Atlantico, Viagens e Turismo SA, Portugal

Co-authors: Helder Miguel Correia Virtuoso Sebastiao, Ana Francisca Ferreira Martins

Two portfolio selection models are compared: the Mean-Variance (MV) model, introduced by Markowitz, and the Mean Absolute Deviation (MAD) model, proposed by Konno and Yamazaki. A large database is used for this comparison, consisting of the 500 largest U.S. stocks, analysed across different time frames. The study period spans from January 4, 2010, to November 30, 2021. After applying rebalancing strategies at various frequencies, an out-of-sample performance analysis is conducted for both models. Several performance metrics are employed in this evaluation, including win rate, cardinality, turnover, annualized mean return, annualized standard deviation, Sharpe ratio, Sortino ratio, Conditional Value at Risk (CVaR), and maximum drawdown. The results indicate that although the minimum-risk solution of the MAD model yields slightly worse values than the corresponding MV model, the differences are relatively small. However, a notable drawback of the MV model is that it tends to allocate capital across all 500 assets, many of which have extremely small weights, and this is an unrealistic assumption in practical investment scenarios. In contrast, the MAD model recommends investing in a more concentrated portfolio of approximately 40-50 assets, which aligns more closely with real-world portfolio management practices.

C1385: **High dimensional hedonic models for spatial price index construction in the accommodation market using web scraped data**

Presenter: **Ilaria Benedetti**, University of Tuscany, Italy

Co-authors: Tiziana Laureti, Niccolo Salvini

The analysis of spatial price disparities is often hampered by the limited availability of detailed information on product characteristics. The aim is to propose a comprehensive framework to overcome this limitation by leveraging high-dimensional webscraped data for the construction of robust spatial price indices (SPI) for differentiated services. The focus is on the Italian accommodation market, where online platforms provide a rich source of granular data. The methodology integrates a scalable web scraping architecture with a multistage data processing pipeline. Data is collected from metasearch engines, capturing a wide array of features for each accommodation, including multiple price offers, customer ratings, and a high-dimensional set of amenities and service attributes. The statistical techniques employed to transform this raw, complex data into an analysis-ready dataset are detailed, including spatial deduplication and robust feature engineering. Building on this dataset, the core of the contribution lies in the application of a high-dimensional hedonic pricing model. Using the detailed scraped characteristics, the model controls for quality differences and isolates price gaps across areas. This yields spatial price indices that are more accurate and stable than those from conventional data. The replicable framework turns unstructured web data into indicators ready for policy analysis and consumer information.

CC418 Room BCB 311 CLUSTERING

Chair: Anahita Nodehi

C1426: **Clustering multivariate discrete data with partial records**

Presenter: **Kevin Giddings**, Carleton University, Canada

Co-authors: Isen McDonald, Utkarsh Dang, Sanjeena Dang

The ability to cluster data with incomplete records is vital in many disciplines. A model-based clustering approach is developed for clustering multivariate discrete data with missing entries using a mixture of multivariate Poisson lognormal distributions. A multivariate Poisson lognormal distribution is a hierarchical Poisson distribution that can account for overdispersion and can model the correlation between the variables. To illustrate the effectiveness of this method, extensive simulation studies are conducted under varying levels of missingness and types of missing data, to evaluate the robustness of this new method under different percentages of incomplete records and missing data patterns. In addition, the approach is performed on a complete zebrafish RNA sequencing dataset. The results obtained from this complete data clustering problem are compared to the performance when some of the count values are artificially omitted. Through this, the method is shown to achieve similar clustering performance between complete and incomplete datasets.

C1422: **Clustering microbiome data using a mixture of logistic normal multinomial distributions**

Presenter: **Sanjeena Dang**, Carleton University, Canada

The human microbiome plays an important role in human health and disease status. Next-generation sequencing technologies allow for quantifying the composition of the human microbiome. Clustering these microbiome data can provide valuable information by identifying underlying patterns across samples. However, clustering these datasets is challenging. Taxa count data in microbiome studies are typically high-dimensional, over-dispersed, and can only reveal relative abundance, and therefore often are treated as compositional. Analyzing such compositional data presents many challenges because they are restricted to a simplex. The aim is to present recent advances in clustering microbiome data using a mixture of logistic normal multinomial models. In a logistic normal multinomial model, the relative abundance of the microbiome is mapped from a simplex to a latent variable in the real Euclidean space using the additive log-ratio transformation. An efficient framework is utilized for parameter estimation using variational Gaussian approximations (VGA). Adopting a variational Gaussian approximation for the posterior of the latent variable reduces the computational overhead substantially. Recent developments using extensions of the LNM distribution to cluster high-dimensional microbiome data are discussed.

C1208: **Incorporating group fairness in a variable-weighted adjacency construction for spectral clustering**

Presenter: **Jesse Ghashti**, University of British Columbia, Canada

Co-authors: John Thompson

Constraining clustering algorithms or applying post-hoc modifications to enforce fairness often pushes the objective function away from local or global optima. A nonparametric method is proposed for group-fair bandwidth selection, motivated by kernel density estimation, where estimated bandwidths act as variable-specific scaling factors in a pairwise weighted distance. These distances are then converted to similarities with a kernel function to form the adjacency matrix for both hard and soft spectral clustering. Applications show that the proposed method produces clustering results that are more group-fair than standard spectral clustering, while yielding smaller deviations from the optimum than other fairness-constrained variants. Finally, it is shown that soft clustering provides probabilistic memberships with a more natural interpretation than hard assignments, and extensions to mixed variable types and alternative fairness definitions are discussed.

C1463: **Towards scalable survival analysis: Efficient clustering via k-means and log-rank**

Presenter: **Nora Villanueva**, University of Vigo, Spain

Co-authors: Marta Sestelo, Luis Machado

Survival analysis provides essential tools for studying time-to-event data, with the comparison of survival curves across groups being one of the main objectives. Traditional clustering approaches often rely on bootstrap-based procedures to approximate the null hypothesis distribution. While effective, they impose heavy computational demands and limit scalability in large datasets. The aim is to present a novel method that leverages k-means and the log-rank test to efficiently identify and cluster survival curves. By eliminating the need for intensive resampling, the approach substantially reduces computation time while preserving statistical validity. Through simulation studies, the proposed method is demonstrated to achieve performance comparable to bootstrap-based clustering techniques, while offering a significant gain in efficiency. These findings highlight that the proposed method offers a practical and scalable alternative for the analysis of multiple survival curves.

CC428 Room BCB 313 COMPLEX AND STRUCTURED TIME SERIES**Chair: Johan Lyhagen****C1259: Higher-order integer autoregression for count time series****Presenter:** Robert Jung, University of Hohenheim, Germany**Co-authors:** Andrew Tremayne

First-order integer-valued autoregressive models are widely regarded as a theoretical workhorse for the analysis of count time series. While the theoretical literature on these models is extensive, relevant applications in empirical work remain relatively scarce. Moreover, little attention has been paid to the consequences of choosing higher-order specifications. The aim is to examine the role of integer-valued autoregressions beyond first order, with particular emphasis on how dynamic properties, both first- and second-order, affect inference and forecasting. Using a combination of real-world data and simulation experiments, the implications of alternative model specifications are assessed. The analysis is facilitated by the newly developed R package *coconots*. The results show that the choice of lag order has non-trivial consequences, and that straightforward extensions of the basic first-order framework are not always appropriate.

C1334: Functional threshold autoregressive model**Presenter:** Yuanbo Li, University of International Business and Economics, China**Co-authors:** Chun Yip Yau

A functional threshold autoregressive model is proposed for flexible functional time series modeling. In particular, the behavior of a function at a given time point can be described by different autoregressive mechanisms depending on the values of a threshold variable at a past time point. Sufficient conditions for the strict stationarity and ergodicity of the functional threshold autoregressive process are investigated. A novel criterion-based method is developed simultaneously, conducting dimension reduction and estimating the thresholds, autoregressive orders, and model parameters. The consistency and asymptotic distributions of the estimators of both thresholds and the underlying autoregressive models are also established. Simulation studies and an application to U.S. Treasury zero-coupon yield rates are provided to illustrate the effectiveness and usefulness of the proposed methodology.

C1279: Mixing conditions for functional autoregressive models and their application to generalization error analysis**Presenter:** Shuntaro Suzuki, Osaka University, Japan**Co-authors:** Yoshikazu Terada

The focus is on nonlinear functional autoregressive models (N-FAR), which are useful for modeling large-scale time series data such as meteorological and temperature fields. First, sufficient conditions are derived under which N-FAR satisfies exponential mixing as its statistical dependence structure, thereby providing a foundation for analyzing the asymptotic behavior of estimators. Next, under these sufficient conditions, we establish variance evaluation when the nonlinear operator is learned by a broad range of estimation methods, including deep learning. Furthermore, as an application, the focus is on the case where the nonlinear operator belongs to the class of integral-type Urysohn operators, and the generalization error is theoretically derived when the integral kernel is estimated using deep learning. The present results provide an integrated guarantee of mixing properties for N-FAR together with variance and generalization error analysis of learners, thereby offering theoretical guidance for practical learning design of functional time series.

C1389: Hierarchical principal component analysis of high-dimensional spatial-temporal data**Presenter:** Chi Tim Ng, Hang Seng University of Hong Kong, Hong Kong

Modern data collection technology allows researchers to obtain multivariate data from different geographical locations and different time points. For example, the Google search frequency of $P = 100$ health-related keywords in $N=52$ states in the United States over a period of $T=52$ weeks. In large-scale tensor-valued datasets, P, T, N can be even greater. Traditional principal component analysis methods of data reduction are designed for the dataset represented as a two-dimensional array, as the rows and columns of a spreadsheet. The extra dimension in the high-dimensional spatial-temporal data (P, T, N) gives new challenges to the researchers. The goal is to propose a hierarchical principal component analysis method for the tensor-valued data of three or even higher-dimensional arrays. By summarizing the huge dataset into relatively few driving forces, the massive dataset can be visualized through the changes in the influences of these driving forces over time and geographical regions. Empirical examples are used to illustrate the proposed hierarchical principal component analysis methods. The performances are evaluated through numerical simulation.

CC384 Room BCB 402 ALGORITHMS AND SOFTWARE**Chair: Jose Perusquia****C1158: A framework for structural variation analysis: Realistic simulation, robust detection, and Haplotype inference****Presenter:** Yongyi Luo, The Chinese University of Hong Kong, China**Co-authors:** Zhen Zhang, Jingyu Hao, Jiandong Shi, Weichuan Yu, Xiaodan Fan

Structural variations (SVs) are critical genomic variants affecting evolution and disease susceptibility. Their analysis faces three major challenges: Existing simulators struggle to capture the complex genomic distribution of SVs; detection methods lack sensitivity for complex structural variants (CSVs) with nested or multi-breakpoint architectures; and haplotype reconstruction suffers from error propagation due to separate variant calling and phasing. To address these challenges, BVSIM is developed, a benchmarking variation simulator that learns empirical distributions of SVs from real genomic data. BVSIM accurately preserves length distributions, telomere-proximal enrichment, and tandem repeat associations seen in human SVs, outperforming existing simulators. Next, gSV is introduced, a general SV detector that combines alignment-based signal decomposition with assembly-based validation. gSV utilizes maximum exact match strategies and graph-cut optimization, enabling sensitive identification of CSVs without pre-defined assumptions about variant types. It excels particularly in detecting nested and multi-breakpoint variants. Finally, DIHap is proposed, a unified probabilistic framework for direct haplotype inference from sequencing data. DIHap simultaneously models sequencing error profiles and haplotype structures, enhancing accuracy in low-coverage and polyploid scenarios.

C1296: FARS: Factor augmented regression scenarios in R**Presenter:** Gian Pietro Enzo Bellocca, Universidad Carlos III de Madrid, Spain**Co-authors:** Ignacio Garron Vedia, Vladimir Rodriguez-Caballero, Esther Ruiz

In the context of macroeconomic/financial time series, the FARS package provides a comprehensive framework in R for the construction of conditional densities of the variable of interest based on the factor-augmented quantile regressions (FA-QRs) methodology, with the factors extracted from multi-level dynamic factor models (ML-DFMs) with potential overlapping group-specific factors. Furthermore, the package also allows the construction of measures of risk as well as modeling and designing economic scenarios based on the conditional densities. In particular, the package enables users to: (i) extract global and group-specific factors using a flexible multi-level factor structure; (ii) compute asymptotically valid confidence regions for the estimated factors, accounting for uncertainty in the factor loadings; (iii) obtain estimates of the parameters of the FA-QRs together with their standard deviations; (iv) recover full predictive conditional densities from quantile forecasts; (v) obtain risk measures based on the extreme quantiles of the conditional densities; and (vi) estimate the conditional density and the corresponding extreme quantiles when the factors are stressed.

C1295: Joint sparse optimization: Methodologies and its applications**Presenter:** Carisa Kwok Wai Yu, The Hang Seng University of Hong Kong, Hong Kong

In many practical scenarios, gathering multiple measurement signals can greatly improve the quality of data analysis. The correlation between these related signals enables the simultaneous analysis of different variables, allowing their influences to be focused on common locations. Joint sparse

optimization (JSO) takes advantage of this principle by capitalizing on the collective effect observed across various measurement signals, thereby enhancing model analysis and sparse recovery capability. The purpose is to explore popular numerical techniques used to tackle JSO challenges, investigating both their theoretical underpinnings and real-world applications. Additionally, a range of numerical tests is performed to assess and compare the efficiency and performance of these techniques, offering valuable insights into their effectiveness across different application areas. Through this investigation, the aim is to deepen the understanding of JSO and its significance in improving data analysis in several disciplines.

C1386: An interactive shiny web application for multivariate analysis of transcriptomic data

Presenter: **Jose Luis Romero Bejar**, University of Granada, Spain

Co-authors: Quintin Mesa Romero, Francisco Javier Esquivel

The growth of transcriptomic technologies such as RNA-seq has produced high-dimensional datasets that require advanced statistical approaches. Multivariate techniques, including principal component analysis (PCA) and cluster analysis, are essential for uncovering biological patterns and reducing complexity. The focus is on reviewing these methodologies and introducing a Shiny-based web application designed to make them accessible to a broader research community. The application enables users to upload gene expression matrices and perform analyses through an intuitive interface, eliminating the need for programming expertise. Core functionalities include data validation, PCA visualization, hierarchical and non-hierarchical clustering with optimal cluster estimation, and differential gene expression analysis with packages such as DESeq2 and edgeR. Interactive plots, tables, and workflow control enhance usability and robustness. This prototype represents a practical step toward democratizing transcriptomic data exploration, bridging the gap between complex computational methods and biological insight in bioinformatics.

CC423 Room BCB 407 NONPARAMETRIC HYPOTHESIS TESTS

Chair: Tiejun Tong

C1151: The Lepage-type statistic for the one-sided location-scale alternative in the two-stage design

Presenter: **Kenji Nitta**, Tokyo University of Science, Japan

Co-authors: Hidetoshi Murakami

In many practical applications, an increase in location is often accompanied by an increase in variance, suggesting that heteroscedasticity may signal a treatment effect. This motivates the use of location-scale tests, which jointly assess differences in both location and scale. A common approach involves combining separate statistics for location and scale into a unified test statistic. In the context of independent Stage 1 and Stage 2 data, two-stage testing procedures are frequently employed. The aim is to extend such procedures to account for potential correlation between the p-values derived from each stage by incorporating first-stage information into the design of the second stage. Through extensive Monte Carlo simulations under a variety of continuous distributions, it is demonstrated that the proposed test statistic achieves enhanced and stabilized power, making it a competitive alternative to existing location-scale test statistics.

C1152: A multisample rank test based on Hermite polynomials

Presenter: **Yuta Sato**, Tokyo University of Science, Japan

Co-authors: Hidetoshi Murakami, Amitava Mukherjee

Multiple aspects involving location, scale, and skewness of statistical distributions should be compared among different samples in practice, instead of only location or scale aspects. While rank statistics based on Legendre polynomials have been studied to this end, a new multisample rank-based test statistic that replaces Legendre polynomials with Hermite polynomials is proposed. The proposed statistic is constructed using the Gram-Schmidt orthonormalization process. The limiting distribution of the proposed test statistic is derived. In addition, simulations are conducted to investigate the convergence of the statistic to its limiting distribution and to compare the power with existing test statistics. Power performances of the proposed test are highly encouraging in various situations compared to some of its competitors. A practical illustration is provided with real data. Some conclusions and future research problems are suggested.

C1460: Testing independence using C-power functions

Presenter: **Mohamed Belalia**, University of Windsor, Canada

Co-authors: Guanjie Lyu

Testing independence within a random vector is a fundamental task in statistical analysis, underpinning numerous applications across scientific domains. With the increasing prominence of copulas in capturing dependence structures, copula-based independence tests have received considerable attention. The purpose is to introduce a novel Cramer-von Mises test statistic based on C-power functions, constructed to detect deviations from independence among components of a continuous random vector. The asymptotic distribution of the test statistic is established under the null hypothesis as well as under a sequence of local alternatives. Evidence from simulation studies and real data applications suggests that the proposed test consistently achieves higher power than the empirical copula-based test, particularly in detecting complex or subtle forms of dependence.

C1364: A new multivariate two-sample rank test based on interpoint distances

Presenter: **Hikaru Yamaguchi**, Tokyo University of Science, Japan

Co-authors: Hidetoshi Murakami

Testing general multivariate two-sample problems is important in many applications. Interpoint distance-based tests are well known to be effective for general alternatives, including simultaneous comparisons of means and variability, or comparisons of the shapes of the distributions. The aim is to propose a new test statistic based on the ranks of interpoint distances to address the two-sample problem for multivariate populations. The proposed test statistic is constructed from the well-known energy distance, and its critical values are approximated using a permutation procedure. Some theoretical properties of the proposed test statistic are also presented. Numerical simulations are conducted to examine the power and robustness of the test under various settings, including high-dimensional data and small sample sizes.

Sunday 14.12.2025

15:50 - 17:30

Parallel Session I – CFE-CMStatistics 2025

CI013 Room BCB 307 EMBEDDINGS IN STATISTICS AND AI**Chair: Joshua Cape****C0179: Counting cycles with DeepSeek***Presenter:* **Tracy Ke**, Harvard University, United States

Despite recent progress, AI still struggles with advanced mathematics. A difficult open problem is considered: How to derive a computationally efficient equivalent form (CEEF) for the cycle count statistic? The CEEF problem does not have known general solutions and requires delicate combinatorics and tedious calculations. Such a task is hard to accomplish by humans, but it is an ideal example where AI can be very helpful. The problem is solved by combining the proposed novel approach and the powerful coding skills of AI. Results use delicate graph theory and contain new formulas for general cases that have not been discovered before. It is found that, while AI is unable to solve the problem all by itself, it is able to solve it if a clear strategy, step-by-step guidance, and carefully written prompts are provided. For simplicity, the focus is on DeepSeek-R1, but other AI approaches can also be investigated.

C0180: How to validate an AI that outperforms humans*Presenter:* **Karl Rohe**, University of Wisconsin-Madison, United States

In various repetitive tasks, modern language models (e.g., ChatGPT) have the potential to exceed the quality of human-generated data. This creates a fundamental challenge in evaluating/validating language models for these tasks: How is a system validated when the validation labels (i.e., from humans) are potentially less accurate than the system's outputs? If humans are the gold standard, then disagreements are errors "blamed on" the AI. A path forward is provided. It is a class of statistical models and algorithms that help better understand and use Cohen's Kappa.

C0181: The manifold hypothesis in science and AI*Presenter:* **Patrick Rubin-Delanchy**, University of Edinburgh, United Kingdom

The manifold hypothesis is a widely accepted tenet of machine learning, which asserts that nominally high-dimensional data are in fact concentrated around a low-dimensional manifold. Some real examples of manifold structure occurring in science and in AI (internal representations of LLMs) are shown, and associated questions are discussed, particularly around how observed topology and geometry might map to the real world (science) or a human-understandable concept (AI). Statistical models and embedding theory are presented, which help explain the efficacy of popular combinations of tools for manifold learning, such as PCA followed by t-SNE. Finally, a vast array of unexplored possibilities in representation learning and potential implications for the future role of AI in science are pointed out.

CO151 Room BCB G07 HiTEC: STATISTICAL LEARNING AND REGIME-SWITCHING IN FINANCE**Chair: Christina Erlwein-Sayer****C1245: Modelling contagious bank runs***Presenter:* **Luitgard Veraart**, London School of Economics and Political Science, United Kingdom

The purpose is to develop a modelling framework for contagion in financial networks arising from bank runs. It is shown how interacting channels of contagion, namely, funding withdrawals in the interbank network and price-mediated contagion arising from fire sales, can turn a bank run on one institution into a systemic crisis. Furthermore, it is also modelled how contagion effects can lead to additional bank runs. The model allows for a wide range of withdrawal mechanisms, both by banks and by external depositors. It can be used for financial stress testing and particularly for analyzing implications of different withdrawal mechanisms for systemic risk. This is illustrated in stylized examples and an empirical case study. It is found that the extent of systemic risk is highly sensitive to the choices of withdrawal strategies used by the market participants. Policy implications are also discussed.

C1317: Valuation of derivatives on carbon emission allowance*Presenter:* **Rogemar Mamon**, University of Western Ontario, Canada

Spot prices of the EU ETS carbon emission allowance are investigated and modelled. A regime-switching mechanism is embedded in various stochastic models as inspired by the volatility clustering phenomenon. The recursive filtering technique is utilized in model calibration. The pricing of European-style futures call options is considered and assessed by comparing pricing errors using the EUA futures data compiled by Bloomberg. A particular hidden Markov model outperforms some benchmarks. The cost-of-carry relationship is also found to hold for both the interphase and intra-phase contracts.

C1375: Risk culture indicators and their role in financial institutions resilience*Presenter:* **Natalie Packham**, Berlin School of Economics and Law, Germany*Co-authors:* Sami Alkhoury

A strong risk culture is generally thought to be valuable to a financial institution as it is said to strengthen an institution's resilience. Can this claim be substantiated? The purpose is to establish quantitative and qualitative risk culture indicators (RCIs) and, using a data set comprising 81 European banks, build a score for risk culture and a stress test score. A relatively better result is found in the 2014 ECB stress test, which corresponds to a better risk culture. Revisiting this problem in the age of generative AI and NLP, it is shown how to build a comprehensive set of risk culture scores from financial institutions' annual reports by training a variety of NLP models. This involves, in particular, reinforcement learning from human feedback (RLHF), which increases the quality of automatically generated RCIs.

C1405: Adaptive forecasting of electricity spot prices: A hybrid HMM-LSTM approach enabling regime detection and prediction*Presenter:* **Christina Erlwein-Sayer**, University of Applied Sciences HTW Berlin, Germany*Co-authors:* Tilman Sayer

Electricity spot prices are characterized by frequent volatility and sudden price jumps, which complicates accurate intraday price forecasting. Forecasting models often struggle to adapt to the dynamic nature of these price movements, as market conditions can shift rapidly. The aim is to propose a novel hybrid approach that combines the predictive power of Long Short-Term Memory (LSTM) networks with the flexibility of Hidden Markov Models (HMMs) to better capture these regime shifts. The LSTM architecture is enhanced by incorporating an underlying HMM that detects and adapts to changing market regimes. These regimes are adaptively filtered from the data, enabling the model to separate the spot price series into regime-specific segments. The recurrent neural network is trained for each regime, and the HMM-derived state probabilities are integrated as gates into the prediction model. The integration of HMM with LSTM forecasts enhances both the accuracy and interpretability of the model. Each LSTM forecast is conditioned on the filtered state of the underlying market regime, which allows the model to capture the different price dynamics that occur in varying market conditions. The model is applied to German spot prices, demonstrating its effectiveness in forecasting price movements by leveraging both short- and long-term market dynamics. The HMM-gated LSTM framework outperforms traditional forecasting methods in terms of predictive accuracy.

CO317 Room BCB G08 ADVANCES IN STRUCTURAL ECONOMETRICS**Chair: Max Antoine Lesellier****C0347: Identifying marginal costs from discrete prices and product characteristics***Presenter:* **Mathieu Marcoux**, Université de Montréal, Canada*Co-authors:* Colin Jaures Ndonfack Zango

A new approach is proposed to identify the marginal costs of differentiated products when prices and other observed product characteristics are discrete. It boils down to a discrete choice problem allowing for flexible distributions of unobservables, therefore avoiding strong distributional assumptions on unobserved demand shocks. The approach treats the conditional expectation of demand as a nuisance function. If data about market shares or sales are available, one can estimate this nuisance function prior to recovering marginal costs. If demand data are not available, a natural exclusion restriction is leveraged between costs and demand, and the nuisance function is approximated using differences of observed product characteristics. The identification procedure features appealing computational properties despite the large number of parameters. Simulation evidence suggests that the ability to precisely recover marginal costs depends on the quality of the approximation of the demand's conditional expectation.

C0564: Sharp estimation of static entry games with covariates

Presenter: **Max Antoine Lesellier**, University of Montreal, Canada

Entry games are a widely used framework for analyzing industry dynamics, using easily accessible data on market entry. New tools are introduced to improve the estimation of static entry games, particularly in scenarios where the equilibrium selection mechanism remains unrestricted. In such settings, the econometrician typically uses inequalities implied by the model to characterize the set of admissible parameters; however, entry games often generate so many inequalities that using them all becomes impractical. To manage this complexity, a recursive algorithm that simultaneously selects a relevant subset of inequalities is developed, and theoretical upper bounds are calculated on the probability of each outcome without relying on simulation-based methods. Additionally, the testing procedure introduced by a prior study is extended, which is consistent and produces an asymptotically pivotal test statistic, to allow for the pre-estimation of the parameters that are point-identified. The approach is particularly useful in handling covariates, including continuous ones. Through comprehensive Monte Carlo simulations, the effectiveness of the estimation procedure is demonstrated. Finally, the method is applied empirically.

C0651: Nonparametric discrete-choice and sensitivity to distributional assumptions on random utility

Presenter: **Jonas Lieber**, Imperial College London, United Kingdom

Co-authors: Alexander Torgovitsky, Pietro Tebaldi

New results are introduced to estimate sharp nonparametric bounds on pre-specified demand counterfactuals in a discrete choice model, avoiding assumptions on the distribution of random utility. The set of inequalities characterizing choice under quasilinear utility points to an efficient linear algebra procedure to partition the space of unobservable, maintaining all relevant information for the question of interest. The method is used to discuss the sensitivity of parametric estimates to the relaxation of widespread distributional assumptions on random utility.

C0735: Identification of dynamic panel logit models with fixed effects

Presenter: **Thomas Russell**, Carleton University, Canada

Co-authors: Christopher Dobronyi, Jiaying Gu, Kyoo il Kim

It is shown that identification in a general class of dynamic panel logit models with fixed effects is related to the truncated moment problem from the mathematics literature. This connection is used to show that the identified set for structural parameters and functionals of the distribution of latent individual effects can be characterized by a finite set of conditional moment equalities subject to a certain set of shape constraints on the model parameters. In addition to providing a general approach to identification, the new characterization can deliver informative bounds in cases where competing methods deliver no identifying restrictions, and can deliver point identification in cases where competing methods deliver partial identification. An estimation and inference procedure is then presented that uses semidefinite programming methods, is applicable with continuous or discrete covariates, and can be used for models that are either point- or partially-identified. Finally, the identification result is illustrated with a number of examples, and an empirical application is provided to employment dynamics using data from the National Longitudinal Survey of Youth.

CO099 Room BCB G09 TOPICS IN MACROFINANCE AND ECONOMETRICS

Chair: Jose Olmo

C0238: Bubbling up? What consumer expectations reveal about U.S. housing market exuberance

Presenter: **Efthymios Pavlidis**, Lancaster University Management School, United Kingdom

Co-authors: Enrique Martinez-Garcia

The presence of speculative bubbles is investigated in the U.S. housing market after the global financial crisis. Unlike standard approaches that rely on observed economic fundamentals, the method leverages subjective price expectations from the University of Michigan Survey of Consumers to test for exuberance without imposing a specific model of intrinsic housing values. By applying recursive least-squares and quantile-based unit root tests to cumulative expectational errors, novel evidence of speculative dynamics is uncovered at the aggregate level and across broad demographic and socioeconomic groups. A date-stamping exercise reveals widespread exuberance in the second half of the 2010s, which paused before the pandemic recession and resurfaced amid the subsequent housing boom in 2021. For the Covid-19 period, notable differences are documented in the timing of exuberance between observed house prices and survey-based indicators- a finding that underscores the importance of controlling for fundamentals when identifying speculative behavior. A complementary analysis using the New York Fed Survey of Consumer Expectations corroborates the baseline results. Overall, findings highlight the value of survey data for monitoring housing markets.

C0266: A robust GRS statistic

Presenter: **William Pouliot**, University of Birmingham, United Kingdom

The GRS statistic is a cross-sectional test of the one-factor capital asset pricing model. That statistic is not appropriate for tests of the K-factor model. Using the same data as prior studies, it is shown that they did not use the specific GRS statistic. To provide clarity on which GRS statistic to use, the detailed mathematical derivation of the cross-sectional variance of the OLS estimators of the estimated intercepts of the K-factor model is provided. This variance is then used to construct the enhanced version of the GRS statistic implemented in the literature. Assumptions underlying that statistic assume invariance of cross-sectional variances and covariances of the idiosyncratic errors in these models. The robust version of the GRS statistic (for the K-factor model) is constructed, allowing for time variation of the variances and covariances of the cross-sectional errors. The consistency of that estimator is established, and then the corresponding asymptotic distribution is also established. When this statistic is calculated on the same data used in a prior study, it allows for a more nuanced comparison between three-, four-, and five-factor models. A comparison of the power functions of the GRS and the GRS developed for the K factor model is also undertaken.

C0343: Shrinkage estimation of spatial panel data models with multiple structural breaks and a multifactor error structure

Presenter: **Chaowen Zheng**, University of Southampton, United Kingdom

Spatial panel data models are investigated with a multifactor error structure and multiple structural breaks occurring in the coefficients of both spatial lagged and explanatory variables. To address the dual challenges of endogeneity and time heterogeneity, a novel penalized generalized method of moments estimation with common correlated effects (PGMM-CCEX) is proposed. Specifically, this method addresses the endogeneity issue by utilizing the cross-sectional averages of regressors as factor proxies when constructing the internal instrumental variables, while employing adaptive group fused lasso to detect multiple structural breaks. The PGMM-CCEX method consistently estimates both the number of breaks and their locations. Furthermore, the post-PGMM-CCEX regime-specific coefficient estimates are consistent and asymptotically follow a normal distribution. Notably, the method remains valid even when factor loadings vary over time, whether synchronously or asynchronously with the parameters of interest. Monte Carlo simulations confirm the satisfactory finite-sample performance of the proposed PGMM-CCEX method. Finally,

the method is applied to analyze cross-country economic growth across 107 countries from 1970 to 2019, revealing the time-varying influence of key economic factors on growth dynamics.

C0823: **Factor-based nowcasting with missing data**

Presenter: **Mark Hallam**, University of York, United Kingdom

Motivated primarily by the problem of producing nowcasts for large unbalanced panels of data, methods are developed for estimating the latent factor structure and values of missing observations in panels with general patterns of missing data, including those with mixed sampling frequencies. The approaches combine recent developments in factor-based imputation for partially observed panels proposed by prior studies with the Kalman filtering and QML-based approaches for estimating approximate dynamic factor models in prior studies. This results in two-step estimators of the latent factors and missing observations that allow for general patterns of missing data, dynamics in the latent factors, and cross-correlation in the idiosyncratic components. Simulation exercises are conducted, using DGPs commonly employed in the factor model literature, modified to allow for various patterns of missing data. It is found that the proposed methods generally perform well, providing accuracy gains relative to the simpler one-step factor-based imputation methods in both the estimation of the latent factors and imputation of the missing values. Finally, the proposed methods are applied to a macroeconomic nowcasting exercise using a real-time US macro-financial dataset over the period 1990 to 2025.

CO269 Room Virtual R01 EMPIRICAL MACRO

Chair: Laura Jackson Young

C0314: **A restricted FAVAR for regional analysis of tariff uncertainty**

Presenter: **Laura Jackson Young**, Bentley University, United States

Co-authors: Michael Owyang

A factor-augmented VAR (FAVAR) is estimated, using a panel of quarterly state- and national-level data to examine how uncertainty related to tariff policy influences state-level economic conditions. The model includes three types of factors: A national factor, regional factors that model the correlation within a particular BEA region, and state-level factors that model the correlation across variables within a particular state. Measures of tariff volatility are included, and the impulse responses of state-level employment, unemployment, hours, and labor force participation are considered.

C0462: **Signaling processing monetary policy surprises**

Presenter: **Sebastian Laumer**, University of Richmond, United States

Co-authors: Italo Santos

High-frequency identification has become the standard approach for identifying monetary policy shocks. Recently, however, this method has come under scrutiny. Several studies show that high-frequency instruments are contaminated by either central bank information effects or the Fed-responds-to-news channel. A methodology is developed to quantify and address both sources of contamination simultaneously. First, it is shown that instruments orthogonalized to one form of contamination remain correlated to the other. Second, a forward regression algorithm is implemented to select the optimal set of predictors. The algorithm consistently chooses a mix of economic news and central bank information variables, suggesting that both channels are present in high-frequency monetary policy surprises. Third, using the selected predictors, distinct shock series are estimated for monetary policy, central bank information effects, and the Fed-responds-to-news channel. It is found that monetary policy shocks reduce output and prices while tightening financial conditions. Fed-responds-to-news shocks, by contrast, raise output and prices and ease financial conditions. Finally, central bank information shocks increase output on impact despite rising interest rates, spreads, and excess bond premia.

C0537: **Corporate tax reforms and the investment channel of monetary policy**

Presenter: **Ezgi Kurt**, Bentley University, United States

Co-authors: Gonzalo Basante

The first empirical evidence on how corporate tax policy affects the effectiveness of monetary policy on investment is presented. By examining exogenous marginal tax reforms in the US, and accounting for the dynamic nature of taxable income and the forward-looking behavior of investment, it is found that monetary policy is more effective at stimulating investment when firms face tax increases, compared to when their taxes are stable.

C0549: **The evolution of community bank interconnectedness**

Presenter: **Giorgi Nikolaishvili**, Wake Forest University, United States

It is found that the community banking sector in the United States has become more interconnected since the global financial crisis, which implies greater exposure to systemic risk and increased vulnerability in future financial crises. A hierarchical dynamic factor model is estimated using a Bayesian approach to extract posterior distributions of national, regional, and state-level latent drivers of quarterly fluctuations in state-average community bank return-on-equity for all 50 US states. The resulting estimates show evidence of both considerable national comovement and state-specific idiosyncrasy with no signs of significant regional comovement. Furthermore, the results show a decrease in the intensity of idiosyncratic dynamics of state-level community bank profitability since the crisis, along with an increase in national comovement across most states.

CO042 Room BCB 206 CFE SESSION: A TRIBUTE TO H. PESARAN II

Chair: James MacKinnon

C0797: **Information matrix tests for switching regression models**

Presenter: **Enrique Sentana**, CEMFI, Spain

Co-authors: Dante Amengual, Gabriele Fiorentini

The EM principle implies that the moments underlying the information matrix test for multivariate switching regression models are the filtered values of the moments that the information matrix test would test if one knew the latent component each observation belongs to. Thus, components are identified related to the conditional heteroskedasticity, conditional and unconditional skewness, and kurtosis of the multivariate regression residuals for each of the regimes. The Monte Carlo simulations indicate that the expressions obtained by numerical integration for the asymptotic covariance matrix of those empirical moments adjusted for sampling variability in the maximum likelihood parameter estimators provide reliable finite sample sizes and good power against various alternatives, especially combined with the parametric bootstrap.

C1371: **An alternative bootstrap procedure for factor-augmented regression models**

Presenter: **Takashi Yamagata**, University of York, United Kingdom

Co-authors: Peiyun Jiang

A novel bootstrap algorithm that is more efficient than existing methods for replicating the asymptotic distribution of the factor-augmented regression estimator for a rotated parameter vector is proposed. The regression is augmented by r factors extracted from a large panel of N variables observed over T time periods. General weak latent factor models with r signal eigenvalues are considered, that may diverge at different rates, N_k^a , where $0 < a_k \leq 1$ for $k = 1, 2, \dots, r$. The asymptotic validity of the bootstrap method is established using not only the conventional data-dependent rotation matrix \hat{H} , but also an alternative data-dependent rotation matrix, \hat{H}_q , which typically exhibits smaller asymptotic bias and achieves a faster convergence rate. Furthermore, the asymptotic validity of the bootstrap is demonstrated under a purely signal-dependent rotation matrix H , which is unique and can be regarded as the population analogue of both \hat{H} and \hat{H}_q . Experimental results provide compelling evidence that the proposed bootstrap method achieves superior performance relative to existing approaches.

C0441: Robust tensor factor analysis in the presence of heavy-tails*Presenter:* **Lorenzo Trapani**, University of Leicester, United Kingdom*Co-authors:* Matteo Barigozzi

The purpose is to consider (robust) inference in the context of a factor model for tensor-valued sequences. The consistency of the estimated common factors and loadings space when using estimators based on minimising quadratic loss functions is studied. Building on the observation that such loss functions are adequate only if sufficiently many moments exist, results are extended to the case of heavy-tailed distributions by considering estimators based on minimising the Huber loss function, which uses an L1-norm weight on outliers. It shows that such class of estimators is robust to the presence of heavy tails, even when only the second moment of the data exists. A modified version of the eigenvalue-ratio principle is also proposed to estimate the dimensions of the core tensor and show the consistency of the resultant estimators without any condition on the relative rates of divergence of the sample size and dimensions. Extensive numerical studies are conducted to show the advantages of the proposed methods over the state-of-the-art ones, especially under the heavy-tailed cases. An import/export dataset of a variety of commodities across multiple countries is analyzed to show the practical usefulness of the proposed robust estimation procedure. An R package RTFA implementing the proposed methods is available on R CRAN.

C0159: A forecast-error taxonomy facing multiple shifts*Presenter:* **David Hendry**, University of Oxford, United Kingdom*Co-authors:* Jennifer L Castle, Jurgen Doornik

Distributional shifts occur frequently, with detrimental impacts on forecast accuracy. Almost all econometric forecasting models are equilibrium correcting, so they are susceptible to systematic forecast failure after equilibrium mean shifts, which could be direct or induced. Unanticipated out-of-sample shifts are a well-known cause of forecast failure, but they later become in-sample shifts and require handling. Previous forecast error taxonomies are extended to include shifts both in-sample and after the forecast origin. The taxonomy reveals which shifts do and do not lead to forecast failure, facing both in-sample and post-forecast origin shifts, also highlighting what features are amenable to rapid correction.

CO132 Room BCB 207 SL FOR FORECASTING, STRUCTURAL REDUCTION, AND REGIME CHANGE**Chair: Moinak Bhaduri****C0452: Envelope matrix autoregressive models***Presenter:* **Tharindu De Alwis**, University of West Florida, United States*Co-authors:* S Yaser Samadi

Matrix-valued data, common in many scientific fields, presents challenges for traditional time series analysis due to overparametrization and the loss of structural information when using vectorization. The matrix autoregressive (MAR) model was developed to overcome these limitations by preserving the original matrix structure, thus reducing dimensions and allowing for clearer data interpretations. However, high-dimensional matrix time series still pose a problem for MAR models because of the large coefficient matrices, making it difficult to discern relevant information. To address this, envelope-based MAR (EMAR) models are proposed. The EMAR approach significantly enhances efficiency in estimation and forecasting by reducing parameters and establishing a connection between the mean function and covariance structure. This is achieved through the use of minimal reducing subspaces of covariance matrices. The asymptotic properties of the estimators are established, and simulation studies (under both normal and non-normal conditions) are conducted to compare their efficiency and accuracy against existing methods. Additionally, the practical effectiveness of the EMAR approach is demonstrated with two real-world applications in economics and business.

C0730: Forecasting and downscaling solar irradiance using transformer models*Presenter:* **Hossein Moradi Rekabdararkolaee**, Bowling Green State University, United States

Accurate forecasting of solar irradiance is essential for integrating solar energy into the grid, particularly in light of the U.S. Federal Energy Regulatory Commission's (FERC) Order 2222, which emphasizes the role of distributed energy resources such as rooftop solar in bulk power system operations. Existing forecasting methods often fail to capture the fine temporal variations in solar irradiance, limiting their effectiveness in informing local solar photovoltaic (PV) installations and optimizing renewable energy integration. A novel approach is presented to solar irradiance downscaling and forecasting using deep learning Transformer models, which are designed to capture complex temporal dependencies and variability. The proposed method uses cubic spline interpolation to downscale solar irradiance data from a 15-minute resolution to a 5-minute local resolution. Subsequently, the Transformer model is applied to forecast the downsampled data. Different Transformer topologies were trained on historical data from Brookings, South Dakota (SD), and tested for a 24-hour forecast based on standard error metrics. The findings highlight the potential of Transformer models for improving solar irradiance forecasting while also shedding light on the challenges and benefits of using raw data for such tasks.

C0748: On detecting changes in certain random intensity-driven point processes through martingale-based repeated testing*Presenter:* **Moinak Bhaduri**, Bentley University, United States

As surfacing, probably eventually, from the pandemic, other crises thwart normalcy: Appalling inequality, climate calamity, the banking crisis, upped possibilities of a fresh Cold War. The enduring motif of the time is incessant chaos. Frequently, much of that chaos results when one type of stationary system gives way to another. Change detection is mainly about estimating these points of deviation. In case a Poisson-type point process carries the system forward, a brand of online detection algorithms is offered, engineered through permutations of trend-switched statistics and a judicious application of false discovery rate control. Certain members of this family that remain asymptotically consistent and close to the ground truth (evidenced through some Hausdorff-similarity) are isolated to pinpoint estimated change locations. Efficient forecasting proves to be a natural corollary. Change point-based clustering tools will also be examined. It is described how such analyses offer concrete definitions to vague objects like COVID waves and measure their enormity.

C1187: TransDFM: A Bayesian transformer-dynamic factor model for stock market analysis*Presenter:* **Dominic Dayta**, Nara Institute of Science and Technology, Japan*Co-authors:* Kazushi Ikeda, Takatomi Kubo

While traditional dynamic factor models (DFM) provide an interpretable framework for decomposing systematic and idiosyncratic behavior in complex systems such as the stock market, they remain restricted by simplifying assumptions (e.g., linearity). Meanwhile, transformer architectures excel at capturing highly non-linear dynamics but often function as "black boxes". The attempt is to combine these two paradigms by introducing the transformer-dynamic factor model (TransDFM), a novel Bayesian framework for financial time series. The DFMs linear factor loading is replaced with a flexible Bayesian neural network embedding and model the latent factor dynamics using a Bayesian structural aligned mixture of VAR (SaMoVAR) process. This leverages the SaMoVAR's structural equivalence to a vector autoregression to enable a principled application of priors over its dynamic attention weights. Applied to the Philippine Stock Exchange (PSE) and S&P 500, the TransDFM demonstrates strong out-of-sample predictive accuracy and generates well-calibrated uncertainty intervals. Furthermore, despite its highly non-linear structure, the model retains economic interpretability: its primary learned factor shows a statistically significant alignment with systematic market risk as measured by the CAPM beta. The TransDFM provides a powerful new tool that combines the predictive performance of deep learning with the structural analysis of classical econometrics.

| | |
|--|------------------------------|
| CO122 Room BCB 208 ADVANCES IN COMPOSITE LIKELIHOOD INFERENCE FOR HIGH-DIMENSIONAL DATA | Chair: Davide Ferrari |
|--|------------------------------|

C0469: Fast and efficient space-time covariance estimation in large datasets by composite likelihood truncation*Presenter:* **Zhendong Huang**, RMIT University, Australia*Co-authors:* Davide Ferrari, Alessandro Casa

Pairwise composite likelihood, a linear combination of pairwise likelihood functions, is a powerful tool for estimating the variogram of space-time random fields. However, the choice of linear coefficients significantly influences the performance of the resulting estimator, both computationally and statistically. The accuracy can deteriorate dramatically when many noisy or highly correlated pairwise scores are included. A new procedure is introduced for selecting and combining pairwise likelihood components from a large set of feasible candidates, while simultaneously estimating the spatial covariance. The proposed method constructs pairwise estimating equations by minimizing an approximate distance from the full likelihood score, subject to a constraint that reflects available computational resources. This results in truncated pairwise estimating equations containing only the most informative partial likelihood score terms. Asymptotic properties of the method are studied, and numerical experiments are performed.

C0653: Pairwise likelihood estimation of mixed effects ordinal data models via stochastic approximations*Presenter:* **Giuseppe Alfonzetti**, University of Udine, Italy*Co-authors:* Ruggero Bellio, Cristiano Varin

A promising inference strategy to deal with the computational challenges posed by mixed effects models for categorical data is composite likelihood. Of particular interest is the case of crossed random effects, where the integrals involved in the likelihood function drastically increase the computational burden when dealing with massive datasets. Standard approaches for composite likelihood estimation, such as pairwise likelihood, substitute the original likelihood with a surrogate one involving a large collection of lower-dimensional integrals. The high number of such integrals typically prevents the scaling of composite likelihood methods on massive applications. The aim is to present a general framework that augments the traditional pairwise likelihood function with a sampling step over the pool of bivariate margins. The proposed procedure can be framed as a stochastic approximation algorithm, which leads to a new scalable estimator based on bivariate margins of the original likelihood function. The proposal is compared to state-of-the-art methods both on synthetic experiments and real data.

C0438: A semiparametric pairwise likelihood for mixed data in copula models*Presenter:* **Gildas Mazo**, INRA, France*Co-authors:* Ekaterina Tomilina, Florence Jaffrezic

Jointly analyzing mixed-type data in high dimensions is a difficult problem, and yet many applications fall into this case. A semiparametric pairwise likelihood method to estimate the parameters in copula models from a sample of independent and identically distributed observations of mixed-type variables is derived and studied theoretically and numerically. Applications in genomics are presented with the Gaussian copula.

C1484: Approximate unbiased sparse estimating functions in high-dimensions*Presenter:* **Giulia Bertagnolli**, Free University of Bozen-Bolzano, Italy*Co-authors:* Alessandro Casa, Davide Ferrari

In many statistical settings, specifying the full likelihood is challenging if not impossible, which has motivated the development of composite likelihood methods. In high-dimensional contexts, there is often an additional need for sparse estimation to reduce model complexity and identify a limited set of significant parameters. Both challenges are addressed within a unified and general framework. Starting from an unbiased estimating function, we solve a penalized minimization problem: The resulting optimal estimating function is sparse, approximately unbiased, and achieves minimal dispersion distance from the classical score function. By imposing sparsity directly on the set of estimating functions, we significantly reduce the number of estimating equations to be solved. We also propose an efficient algorithm for solving the minimization problem, with two main computational advantages: it operates in a block-wise fashion and avoids inverting large covariance matrices.

| |
|--|
| CO066 Room BCB 209 DIRECTIONAL STATISTICS |
|--|

Chair: Anahita Nodehi**C0675: On regime changes in text data using hidden Markov model of contaminated vMF distribution***Presenter:* **Yingying Zhang**, Western Michigan University, United States*Co-authors:* Shuchismita Sarkar, Yuanyuan Chen, Xuwen Zhu

The aim is to present a novel methodology for analyzing temporal directional data with scatter and heavy tails. A hidden Markov model with contaminated von Mises-Fisher emission distribution is developed. The model is implemented using a forward and backward selection approach that provides additional flexibility for contaminated as well as non-contaminated data. The utility of the method for finding homogeneous time blocks is demonstrated on several experimental settings and two real-life text data sets containing presidential addresses and corporate financial statements, respectively.

C0683: Bayesian model selection for analyzing predictor-dependent directional data*Presenter:* **Ingrid Guevara**, Pontificia Universidad Católica de Chile, Chile*Co-authors:* Vanda Inacio, Luis Gutierrez

Developing statistical models for directional data is increasingly important given the growing need to analyze peak hours in 24-hour services, such as Seoul's bike rental system. Motivated by the aim of identifying factors that influence demand across distinct hours, a model is introduced based on a linear dependent Dirichlet process mixture of projected normal distributions. The approach flexibly captures asymmetric and multimodal densities, while incorporating discrete spike-and-slab priors to facilitate model selection and allow model averaging to account for model uncertainty. The simulation study shows that, across various scenarios, our model successfully recovers the true functional form of the conditional density while reliably selecting the correct model structure, improving accuracy as the sample size increases. The application of the method to the data from the Seoul bike-sharing system successfully unveils that weather conditions significantly impact bike demand fluctuations across distinct hours. The approach also allows predicting peak rental times, revealing that, for instance, on a typical summer day, bike demand decreases between 8 am and 4 pm, while in winter, it drops during the early morning.

C0927: Conditional von Mises Bayesian networks*Presenter:* **Agnese Panzera**, Università degli Studi di Firenze, Italy*Co-authors:* Anna Gottard

Directed acyclic graphical models, commonly known as Bayesian networks, use directed acyclic graphs to represent conditional independence among random variables. A new class of Bayesian networks is proposed for circular random variables, based on the von Mises distribution. The approach is illustrated by applying these models to study the conditional independencies in sequences of angles describing protein structures.

C1018: A robust estimator of location on spheres and manifolds*Presenter:* **Jongmin Lee**, Pusan National University, Korea, South*Co-authors:* Sungkyu Jung

The aim is to introduce Huber means on Riemannian manifolds, including circles and spheres, providing a robust alternative to the Frechet mean by integrating elements of both squared and absolute loss functions. The Huber means are designed to be highly resistant to outliers while maintaining efficiency, making it a valuable generalization of Huber's M-estimator for manifold-valued data. The statistical and computational aspects of Huber

means are comprehensively investigated, demonstrating their utility in manifold-valued data analysis. Specifically, the Huber means are consistent and enjoy the central limit theorem. Additionally, a novel moment-based estimator is proposed for the limiting covariance matrix, which is used to construct a robust one-sample location test procedure and an approximate confidence region for location parameters. The Huber mean is shown to be highly robust and more efficient than the Frechet mean in the presence of outliers or under heavy-tailed distributions on spheres.

CO093 Room BCB 210 TIME SERIES AND STRUCTURAL BREAKS
Chair: Clifford Lam
C0681: Monitoring for a phase transition in a time series of Wigner matrices

Presenter: **Nina Doernemann**, Aarhus University, Denmark

Co-authors: Tim Kutta, Piotr Kokoszka, Sunmin Lee

Methodology and theory are developed for the detection of a phase transition in a time series of high-dimensional random matrices. In the model studied, at each time point $t = 1, 2, \dots$, a deformed Wigner matrix \mathbf{M}_t is observed, where the unobservable deformation represents a latent signal. This signal is detectable only in the supercritical regime, and the objective is to detect the transition to this regime in real time, as new matrix-valued observations arrive. The approach is based on a partial sum process of extremal eigenvalues of \mathbf{M}_t , and its theoretical analysis combines state-of-the-art tools from random-matrix-theory and Gaussian approximations. The resulting detector is self-normalized, which ensures appropriate scaling for convergence and a pivotal limit, without any additional parameter estimation. Simulations show excellent performance for varying dimensions. Applications to pollution monitoring and social interactions in primates illustrate the usefulness of the approach.

C0852: Monitoring time series with short detection delay

Presenter: **Tim Kutta**, Aarhus University, Denmark

Co-authors: Nina Doernemann

Sequential tests are discussed for a change in the mean of a dependent, Banach space-valued time series. For this purpose, a new class of weighted CUSUM statistics is introduced, which is tailored to the detection of changes shortly after they occur. Unlike current alternatives, which either experience long detection delays or offer short delays only at the very beginning of the monitoring period, the presented approach provides consistently short detection delays anywhere in the monitoring period. This property is highly relevant for modern applications, such as epidemiology and finance, where short delays are crucial and the timing of the change is unpredictable. The theoretical results are based on new Hoelderian invariance principles that are proven under some high-level conditions for Banach space-valued data.

C1063: The impact of climate on respiratory drug demand in Greece: A multi-method modelling approach

Presenter: **Matteo Farne**, University of Bologna, Italy

Co-authors: Viviana Schisa

Climate change poses a growing challenge to public health, yet its impact on pharmaceutical demand remains largely overlooked. The aim is to examine the relationship between climate variability and the weekly consumption of prescription respiratory drugs in Greece. Using high-frequency panel data from 20 regions (2016-2023), a multi-method approach is adopted. Structural breaks are first identified in drug demand, identifying a sharp decline during the COVID-19 pandemic, followed by a gradual post-pandemic recovery. Causal relationships are then investigated between climate variables and drug use through spectral Granger causality, uncovering frequency-specific associations that inform model selection and forecasting design. Building on these insights, several forecasting models are developed and compared, including VARX, random forests with moving block bootstrap, and LSTM neural networks. To complement this, a fixed effects lag-distributed regression model that captures the dynamic response is estimated to climate variability while accounting for unobserved spatial and temporal heterogeneity. A spatial lag specification is also employed to detect interregional spillovers in pharmaceutical consumption. Results show that climate variables significantly enhance prediction accuracy, underscoring their value for pharmaceutical planning. Anticipating climate-sensitive demand is key to strengthening healthcare system resilience under evolving environmental pressures.

C1094: High-dimensional sequential change detection

Presenter: **Dogyoon Song**, University of California Davis, United States

The aim is to address the quickest change detection (QCD) problem for multivariate Gaussian time series in which post-change parameters are unknown and must be estimated from a limited data window. The focus is on the high-dimensional regime where the dimension p grows proportionally with the window size n . Extending the window-limited CuSum (WLCuSum) procedure to this setting, it is shown that its asymptotic performance is governed by a novel information-theoretic metric, which is termed the normalized high-dimensional Kullback-Leibler (NHDKL) divergence. Specifically, the detection delay is inversely proportional to the difference between the NHDKL of the post-change versus pre-change distributions and the NHDKL attributable to estimation error. This characterization reveals the suboptimality of plug-in estimators. By minimizing the estimation-error component of the NHDKL, the Ledoit-Wolf quadratic inverse Steins shrinkage estimator (LWISE) is identified as asymptotically optimal within a broad class of shrinkage estimators. Coupled with the sample mean, LWISE yields a practical WLCuSum detector that provably achieves the optimal tradeoff between detection delay and false-alarm rate.

CO022 Room BCB 211 LEARNING AND DISCOVERY IN HIGH-DIMENSIONAL AND STRUCTURED DATA
Chair: Doudou Zhou
C0986: MATES: Multi-view aggregated two-sample test

Presenter: **Doudou Zhou**, National University of Singapore, Singapore

Co-authors: Zexi Cai

The two-sample test is a fundamental problem in statistics with a wide range of applications. In the realm of high-dimensional data, nonparametric methods have gained prominence due to their flexibility and minimal distributional assumptions. However, many existing methods tend to be more effective when the two distributions differ primarily in their first and/or second moments. In many real-world scenarios, distributional differences may arise in higher-order moments, rendering traditional methods less powerful. To address this limitation, a novel framework is proposed to aggregate information from multiple moments to build a test statistic. Each moment is regarded as one view of the data and contributes to the detection of some specific type of discrepancy, thus allowing the test statistic to capture more complex distributional differences. The novel multi-view aggregated two-sample TESt (MATES) leverages a graph-based approach, where the test statistic is constructed from the weighted similarity graphs of the pooled sample. Under mild conditions on the multi-view weighted similarity graphs, theoretical properties of MATES are established, including a distribution-free limiting distribution under the null hypothesis, which enables straightforward type-I error control.

C1031: Robust sensitivity analysis for inverse probability weighting estimation via augmented percentile bootstrap

Presenter: **Xinran Li**, University of Chicago, United States

The identification of causal effects in observational studies typically relies on two standard assumptions: unconfoundedness and overlap. However, both assumptions are often questionable in practice: Unconfoundedness is inherently untestable, and overlap may fail in the presence of extreme unmeasured confounding. While various approaches have been developed to address unmeasured confounding and extreme propensity scores separately, few methods accommodate simultaneous violations of both assumptions. The aim is to propose a sensitivity analysis framework that relaxes both unconfoundedness and overlap, building upon the marginal sensitivity model. Specifically, the bound is allowed on unmeasured confounding to hold for only a subset of the population, thereby accommodating heterogeneity in confounding and allowing treatment probabilities to be zero or one. Moreover, unlike prior work, the approach does not require bounded outcomes and focuses on overlap-weighted average treatment effects, which are both practically meaningful and robust to non-overlap. Computationally efficient methods to obtain worst-case bounds

are developed via linear programming, and a novel augmented percentile bootstrap procedure is introduced for statistical inference. This bootstrap method handles parameters defined through over-identified estimating equations involving unobserved variables and may be of independent interest.

C1073: **Transfer learning for high-dimensional reduced rank time series models**

Presenter: **Abolfazl Safikhani**, George Mason University, United States

The objective of transfer learning is to enhance estimation and inference in a target data by leveraging knowledge gained from additional sources. Recent studies have explored transfer learning for independent observations in complex, high-dimensional models assuming sparsity, yet research on time series models remains limited. The focus is on transfer learning for sequences of observations with temporal dependencies and a more intricate model parameter structure. Specifically, the vector autoregressive model (VAR) is investigated, a widely recognized model for time series data, where the transition matrix can be deconstructed into a combination of a sparse matrix and a low-rank one. The aim is to propose a new transfer learning algorithm tailored for estimating high-dimensional VAR models characterized by low-rank and sparse structures. Additionally, a novel approach is presented for selecting informative observations from auxiliary datasets. Theoretical guarantees are established, encompassing model parameter consistency, informative set selection, and the asymptotic distribution of estimators under mild conditions. The latter facilitates the construction of entry-wise confidence intervals for model parameters. Finally, the empirical efficacy of the methodologies is demonstrated through both simulated and real-world datasets.

C1287: **Monoculture and social welfare of the algorithmic personalization market under competition**

Presenter: **Xin Tong**, University of Hong Kong, Hong Kong

Algorithmic personalization markets, where providers utilize algorithms to predict user types and personalize products or services, have become increasingly prevalent in our daily lives. The adoption of more accurate algorithms holds the promise of improving social welfare through enhanced predictive accuracy. However, concerns have been raised about algorithmic monoculture, where all providers adopt the same algorithms. The prevalence of a single algorithm can hinder social welfare due to the resulting homogeneity of available products or services. The purpose is to address the emergence of algorithmic monoculture from the perspective of providers' behavior under competition in the algorithmic personalization market. It is found that competition among providers could mitigate monoculture, thereby enhancing social welfare in the algorithmic personalization market. By examining the impact of competition on algorithmic diversity, the contribution is to a deeper understanding of the dynamics within algorithmic personalization markets and offers insights into strategies for promoting social welfare in these contexts.

CO289 Room BCB 212 TREATMENT EFFECTS AND POLICY EVALUATION

Chair: Juan Manuel Rodriguez-Poo

C0334: **Policy evaluation with many outcomes**

Presenter: **Maria Nareklshvili**, Stanford University, United States

Probabilistic causal inference is proposed for observational data. The test statistics are designed to discover treatment effects in the presence of heterogeneous effects, noisy outcomes, heavy tails, or complex nonlinear relationships among outcomes. The key idea is to predict treatment status using all outcome variables jointly or individually under an appropriate balancing mechanism or residualization scheme, and test whether the distribution of predicted treatment scores differs significantly between treated and control groups. Simulations and empirical results demonstrate that the method is highly powerful and outperforms the Wald test in settings with small samples, heterogeneous treatment effects, or extreme values.

C0844: **A signature synthetic control method for unbalanced panels**

Presenter: **Pietro Emilio Spini**, University of Bristol, United Kingdom

Co-authors: Stefan Hubner

The aim is to generalize the synthetic control methods to unbalanced panels where the overlap in the time index sets is very limited or absent altogether. Empirically relevant cases include control units with missing data and systematic differences in sample periods and/or frequencies. The approach relies on signatures, a nonparametric feature extraction method. After extracting the signatures of treatment and potential control units in the pre-intervention period, a synthetic control is constructed directly on the signatures. It is shown that the generalized method collapses to standard synthetic controls on a balanced panel. The properties of the resulting synthetic-control estimator are studied for the average treatment effect on the treated, and its convergence rates under standard conditions on the nonparametric component, which is treated as a nuisance function.

C0857: **Testing sign agreement**

Presenter: **Deborah Kim**, University of Warwick, United Kingdom

The focus is on the problem of testing sign agreement among a finite number of means. This problem naturally arises in numerous empirical contexts, including detecting sign-opposing average treatment effects across subgroups in randomised controlled trials, testing the homogeneity of local average treatment effects (LATE) across subpopulations, and examining the testable implications of the assumptions underlying LATE. For the null hypothesis that the means are all non-negative or all non-positive, two novel bootstrap tests are proposed: The least favourable and the hybrid tests. Unlike existing procedures, both tests accommodate arbitrary dependence among estimators for any finite number of means. It is shown that both tests control their asymptotic sizes uniformly over a large class of nonparametric distributions. Results from simulation studies in finite samples indicate that the rejection probabilities of both tests attain the nominal level under the null. The Hybrid test exhibits higher power than the least favorable test against a broad class of alternative distributions, though the opposite holds when comparing only two means. An empirical application to disapproval ratings of the Trump administration demonstrates the practical utility of both tests.

C0923: **Endogeneity in panel and network models: Identification without IV**

Presenter: **Ben Deaner**, UCL, United Kingdom

Co-authors: Andrei Zeleneev

The aim is to present new identification results and estimation methods for causal effects in large N and T panel data.

CO294 Room BCB M201 ADVANCES IN EXTREME VALUE ANALYSIS

Chair: Carlotta Pacifici

C0348: **On predicting the likelihood of high-frequency extreme price movements**

Presenter: **Julien Hambuckers**, University of Liege, Belgium

Co-authors: Philippe Hubner

In finance, the ability to better understand and predict high-frequency extreme price movements (EPM) is crucial for market stability and investor confidence. This is, however, a challenging task, due to the complexity surrounding EPM identification and the numerous predictors available. Extreme value theory is exploited to bypass some of these issues: Rather than focusing on the EPM, the time-varying characteristics and determinants of their distribution are studied. To do so, it is assumed that the distribution for the maxima of high-frequency negative returns over small time intervals can well approximate it. Thus, these maxima are specified for: A parametric generalized autoregressive score model with a generalized extreme value response distribution. To perform covariate selection, a L0-penalized likelihood procedure is introduced. The properties of the estimation and selection procedures are investigated in an extensive Monte Carlo simulation, highlighting good finite-sample performance. Empirically, limit order book data is used from 2009 to 2019 for 50 NASDAQ-traded stocks and a large set of liquidity measures. It is shown that the model performs better than other benchmark specifications, mainly due to the incorporation of important conditioning variables. From the

variable selection procedure, it is shown that liquidity measures play a significant role in determining the distribution of future EPM, along with daily uncertainty indicators.

C0360: **Extremes of structural causal models**

Presenter: **Nicola Gnecco**, Imperial College London, United Kingdom

Co-authors: Sebastian Engelke, Frank Roettger

The behavior of extreme observations is well-understood for time series or spatial data, but little is known about whether the data-generating process is a structural causal model (SCM). The behavior of extremes is studied in this model class, both for the observational distribution and under extremal interventions. It is shown that under suitable regularity conditions on the structure functions, the extremal behavior is described by a multivariate Pareto distribution, which can be represented as a new SCM on an extremal graph. Importantly, the latter is a sub-graph of the graph in the original SCM, which means that causal links can disappear in the tails. A directed version of extremal graphical models is further introduced, and it is shown that an extremal SCM satisfies the corresponding Markov properties. Based on a new test of extremal conditional independence, two algorithms are proposed for learning the extremal causal structure from data. The first is an extremal version of the PC-algorithm, and the second is a pruning algorithm that removes edges from the original graph to consistently recover the extremal graph. The methods are illustrated on river data with known causal ground truth.

C0674: **Accounting for missing data when modelling block maxima**

Presenter: **Emma Simpson**, University College London, United Kingdom

Co-authors: Paul Northrop

Modeling block maxima using the generalized extreme value (GEV) distribution is a classical and widely used method for studying univariate extremes. It allows for theoretically motivated estimation of return levels, including extrapolation beyond the range of observed data. A frequently overlooked challenge in applying this methodology comes from handling datasets containing missing values. In this case, one cannot be sure whether the true maximum has been recorded in each block, and simply ignoring the issue can lead to biased parameter estimators and underestimated return levels. An extension of the standard block maxima approach is proposed to overcome such missing data issues. This is achieved by explicitly accounting for the proportion of missing values in each block within the GEV model. Inference is carried out using likelihood-based techniques, and an update is proposed to commonly used diagnostic plots to assess model fit. The performance of the method is assessed via a simulation study, with results that are competitive with the "ideal" case of having no missing values. The methodology is further demonstrated on ozone data from Plymouth, U.K.

C1440: **On the selection of the threshold with probability weighted moments**

Presenter: **Frederico Caeiro**, NOVA.ID.FCT - Universidade Nova de Lisboa, Portugal

Co-authors: Ivette Gomes

In hydrology and other applied fields, probability-weighted moments (PWMs) are a common tool for estimating the parameters of extreme value distributions. The focus is on the PWM estimator of the extreme value index (EVI) of a Pareto-type model. However, due to the estimator's specific properties, a direct estimation of the threshold is not straightforward. An adaptive approach is studied to choose the number of order statistics to be used in the estimation. Furthermore, the introduced methodology is applied to a dataset in the insurance field.

CO082 Room BCB M202 SUSTAINABLE AND CLIMATE FINANCE

Chair: Michele Costola

C0595: **How climate risk shapes U.S. systemic stability through syndicated lending**

Presenter: **Ana Sina**, University of Reading, United Kingdom

Co-authors: Monica Billio, Alfonso Dufour

The purpose is to explore how climate risk shapes systemic stability in the U.S. syndicated loan market through a framework integrating 25 climate-related indicators adjusted for banks' geographic lending exposures. It combines eight established climate risk measures with a novel proxy for banks' green lending orientation and 16 region-specific indicators capturing anomalies, droughts, and extreme events weighted by loan distribution. Among this wide set of variables, four key factors drive systemic risk: Banks' environmental stance, extreme events in the South, precipitation anomalies in the East, and changes in U.S. climate policy. Systemic risk rises with exposure to climate-vulnerable regions but is mitigated by green lending. While climate policies are vital, they may cause short-term instability, stressing the need for a credible, rule-based transition. The findings call for independent regulation to embed climate risk in finance beyond political fluctuations.

C1017: **Why the ESG boom?**

Presenter: **Marcella Lucchetta**, university of Venice, Italy

The global ESG (Environmental, Social, Governance) investment market surged to \$28.36 trillion in 2024, driven by investor demand and regulatory momentum. The purpose is to investigate the role of ambiguity aversion in ESG asset allocation under uncertain regulatory frameworks, proposing that ambiguity, reflecting market sentiment, shapes investment decisions beyond traditional risk models. A theoretical model of decision-making is developed between ESG and traditional assets, predicting an ESG gap due to underinvestment under ambiguity. Contrary to expectations, a random-effects GLS regression on seven major U.S. ESG equity funds from January 2018 to December 2024 (581 monthly observations) reveals that ambiguity, proxied by Baker's news-based EPU index, boosts ESG returns by 1.62 percentage points per 0.1 EPU increase. Market risk (VIX) reduces returns, while SEC's 2021 disclosure requirements unexpectedly depress performance, suggesting compliance costs. Excess returns over the S&P 500 explain 71.7% of return variation, positioning ESG as a safe-haven asset. This reframes ambiguity as a catalyst for the ESG boom, challenging risk-centric models and offering implications for financial theory and regulatory design. Policymakers should balance transparency with flexibility to support ESG growth, while investors can leverage ESG's resilience in uncertain markets.

C1221: **Measuring climate alignment: Disagreement across metrics and its effect on the cost of capital**

Presenter: **Carmelo Latino**, Smith School of Enterprise and the Environment, University of Oxford and Leibniz Institute for Financial Research

SAFE, Germany

Co-authors: Gireesh Shrimali

The purpose is to review the methodologies of five prominent Paris Alignment Metrics (PAM) and document substantial disagreement across providers. To structure these differences, a taxonomy of methodological choices is proposed, and the main sources of divergence are decomposed. Applying this framework, it is shown that PAMs yield materially different alignment estimates for the same firms. Their financial relevance is then examined by linking PAM scores and their divergence to firms' cost of capital and institutional ownership. The results indicate that stronger and more consistent alignment is associated with a lower cost of capital and greater institutional holdings, while disagreement among different PAMs weakens these relationships.

C1303: **From rain to ruin? Flood impacts on Italian SME loans**

Presenter: **Michele Costola**, Ca' Foscari University of Venice, Italy

Co-authors: Ilaria Prosdocimi, Giacomo Sarrocco

Floods are the most frequent natural hazard in Europe, and their intensification under climate change poses growing risks to firms and financial stability. Small and medium-sized enterprises (SMEs), the backbone of the Italian economy, are especially exposed. The purpose is to investigate how floods affect SME loan performance in Italy by linking loan-level data on 24,219 securitized loans (2004-2023) with a dynamic measure

of flood exposure at the provincial level. Unlike static hazard maps, this approach measures realized physical risk during the life of each loan, capturing the actual timing and severity of flood events. Using survival models, it is found that months of severe flooding increase SME default risk by up to 45 percent, with small firms and sectors such as real estate and water supply particularly vulnerable. By showing how realized climate shocks translate into loan defaults, the results identify a direct channel through which physical risks threaten financial stability and inform stress-testing exercises.

CO119 Room BCB 308 THE PRICE OF INSIGHT: COST-AWARE AND HUMAN-GUIDED LEARNING
Chair: Marshall Honaker
C0426: Learning to defer to multiple experts: From individuals to populations and crowds

Presenter: **Eric Nalisnick**, Johns Hopkins University, United States

Artificial intelligence is being deployed in ever more consequential settings such as healthcare and autonomous driving. Thus, it is ensured that these systems are safe and trustworthy. One near-term solution is to ensure that a human is involved in the decision-making process and that the system can ask for help in difficult or high-risk scenarios. Recent advances are presented in the "learning to defer" paradigm: Decision-making responsibility is allocated to either a human or a model, depending on which is more likely to take the correction action. In particular, novel formulations are presented that can support multiple human decision makers, which could range from known individuals to anonymous members of a crowd.

C0480: Learning to partially defer for sequences

Presenter: **Sahana Rayan**, University of Michigan, United States

In the learning to defer (L2D) framework, a prediction model can either make a prediction or defer it to an expert, as determined by a rejector. Current L2D methods train the rejector to decide whether to reject the entire prediction, which is not desirable when the model predicts long sequences. An L2D setting is presented for sequence outputs where the system can defer specific outputs of the whole model prediction to an expert in an effort to interleave the expert and machine throughout the prediction. Two types of model-based post-hoc rejectors are proposed for pre-trained predictors: A token-level rejector, which defers specific token predictions to experts with next token prediction capabilities, and a one-time rejector for experts without such abilities, which defers the remaining sequence from a specific point onward. In the experiments, it is also empirically demonstrated that such granular deferrals achieve better cost-accuracy tradeoffs than whole deferrals on traveling salesman solvers and news summarization models.

C0514: Iterative modeling under feature acquisition constraints

Presenter: **Marshall Honaker**, University of Pittsburgh, United States

Co-authors: Lucas Mentch, Meredith Wallace

The traditional approach to statistical modeling involves collecting a set of relevant features, performing some exploratory data analysis, and fitting a model to the response. Once fit, the model is used to evaluate new observations and obtain predictions. An implicit but often overlooked assumption in this setup is that all features must be collected for every new observation. In practice, the real-world cost or burden of acquiring a feature's value must be balanced with its utility as a predictor. For example, medical professionals may employ different tests in different situations based on their availability, affordability, and invasiveness. A novel perspective on cost-sensitive learning is proposed that sequentially introduces features in order of increasing cost. Rather than a single model that balances accuracy and feature cost, a sequence of models is produced that optimizes performance subject to the allotted cost at each stage. Moreover, each model is equipped with a reject option so that predictions are made when model confidence is sufficiently high and deferred to the next most costly model otherwise. In medical contexts, this allows patients to be diagnosed using the least costly measures without sacrificing predictive accuracy. Numerous demonstrations of this model-agnostic framework are provided.

C0606: Human-AI collaboration with partial feedback

Presenter: **Ruijiang Gao**, University of Texas at Dallas, United States

As AI becomes increasingly embedded in real-world applications, human-AI collaboration is a daily reality for many workers across diverse fields. Purely autonomous AI systems are rare in practice, making it essential to develop AI systems that can work seamlessly and effectively with human partners. Despite this, much existing research focuses on human-AI collaboration in settings with idealized, full-information feedback, such as access to complete classification labels. In contrast, many practical scenarios only provide partial feedback, where outcomes are observed only for chosen actions, complicating the assessment of both human and algorithmic performance. This limitation is common in business settings like customer service, loan applications, and healthcare. The recent work is discussed on addressing key challenges where humans can help improve traditional algorithmic policy learning under partial feedback: learning from offline observational data, and online contextual bandit problems, where human decision-makers possess valuable but imperfect domain knowledge about optimal actions. By addressing these obstacles, the proposed research seeks to advance responsible AI deployment, enabling more effective human-AI systems that are robust, ethically aligned, and practical for complex real-world applications.

CO251 Room BCB 309 MARKOV-SWITCHING MODELS
Chair: Francesca Loria
C0450: Estimating state-dependent hysteresis effects

Presenter: **Heejee Chang**, University of Washington, United States

Co-authors: Chang-Jin Kim

The potential asymmetry in hysteresis effects is explored between recessions and booms in the U.S. economy. This issue is examined by explicitly modeling the asymmetry within a bivariate VAR framework with Markov-switching volatility. The model allows for the possible correlation between aggregate demand and supply shocks, which captures the degree of hysteresis effects. Identification is achieved through a combination of short-run and long-run restrictions, along with a heteroskedasticity-based approach. Empirical results indicate no evidence of hysteresis effects during either recessions or booms prior to the mid-1980s. During this period, only aggregate supply shocks had permanent effects on output, and recessions were primarily driven by transitory demand shocks. In contrast, since the mid-1980s, clear evidence of asymmetry is found: Hysteresis effects are strong during recessions, but not during booms. In this later period, permanent demand shocks are the main driver of output fluctuations in recessions, while supply shocks dominate in expansions.

C0820: Testing growth vulnerabilities and macro-financial linkages

Presenter: **Leonardo Iania**, UCLouvain, Belgium

Co-authors: Francesco Furno, Francesca Loria, Christian Schulz, Domenico Giannone

The existence of Macro-financial linkages is assessed by testing if financial conditions help predicting macroeconomic risk. The analysis focuses on (predictive) quantile regressions of real GDP on financial conditions, which we reassess by using new robust methods of inference for quantile regression in a time series setting.

C0845: Forecasting with time-varying order-invariant structural vector autoregressions

Presenter: **Tomasz Wozniak**, University of Melbourne, Australia

Co-authors: Annika Camehl

Recent developments suggest that heteroskedastic structural vector autoregressions forecast better when their contemporaneous effect matrix is

time-varying or invariant to the ordering of the variables in the system. However, combining these two features is challenging because order-invariant specifications can be estimated when the model is identified by time-varying volatility or non-normal shocks, a feature that is difficult to ensure in time-varying models. A new forecasting model is proposed that combines fast-moving stochastic volatility with a Markov-switching structural matrix following persistent regimes. This model enables the identification of structural matrices through heteroskedasticity and non-normality within each regime. Additional flexibility is provided by estimating an overfitting number of regimes and extensive hierarchical prior structures. It is demonstrated that both features, time variation and order invariance of the structural matrix, contribute to improvements in the forecasting precision of macroeconomic systems with up to 20 variables.

C0922: Challenges in achieving explainability and control with supply chain forecasts

Presenter: **Rishab Guha**, Amazon, United States

Co-authors: Andrea Tambalotti, Phillip Jang

AI and deep learning methods have revolutionized many forecasting applications but have not achieved widespread adoption in industry for "macro" forecasting (e.g., forecasting aggregate revenue). This paper identifies three critical capabilities that traditional macroeconometrics methods achieve but current AI approaches lack: (1) multivariate consistency at scale, (2) explainable and controllable long-run assumptions, and (3) flexible incorporation of forward-looking external inputs. A Bayesian vector autoregression state-space framework is described, which builds on models used in macroeconomic forecasting, and is used in production at a major e-commerce retailer, where the forecasts influence billions of dollars in spending decisions. By detailing how traditional time series methods solve these challenges today, concrete opportunities for researchers are identified to develop hybrid approaches that combine the accuracy advantages of modern AI with the explainability and control benefits of traditional methods.

CO124 Room BCB 310 RECENT ADVANCES IN LEARNING AND COMPLEX BIOMEDICAL DATA

Chair: Sarbojit Roy

C1204: Causal estimation of cluster-specific effects on spatially associated survival data using SoftBART

Presenter: **Debajyoti Sinha**, Florida State University, United States

The aim is to propose a novel Bayesian approach to estimate causal effects in spatially clustered survival data. Using soft Bayesian additive regression trees (SBART), a nonparametric regression is introduced for a log-normal survival model that accommodates spatial associations among unknown cluster effects through a directed acyclic graph autoregressive (DAGAR) model. A two-stage approach is employed, which entails estimating the propensity score in the first step and incorporating it as a confounder in the outcome model in the second step. The simulation study compares the method with existing approaches under various simulation scenarios, including correctly specified as well as misspecified outcome models, to demonstrate the satisfactory performance of the method. The method is applied to analyze the causal effect of treatment delay (TD) on post-treatment survival of breast cancer patients from the Florida Cancer Registry (FCR). The analysis produces the county-specific as well as state-wide assessment of the causal effects while accommodating spatial association among counties.

C1494: Joint modeling of evoked and induced event-related spectral perturbations

Presenter: **Damla Senturk**, University of California Los Angeles, United States

Event-related spectral perturbations (ERSPs) capture dynamic changes in electroencephalography (EEG) power across frequency and trial time. Even though they are obtained at the trial level, they are commonly averaged across trials and analyzed at the subject level for enhancing the signal-to-noise ratio. While evoked activity is stimulus-locked, representing the brain's predictable response to stimuli, induced signals that are not strictly locked to stimulus presentation are thought to be generated by higher-order processes, such as attention and integration. Motivated by joint modeling of multilevel (trials nested in subjects) and multivariate (evoked and induced) ERSP data from a visual-evoked potentials (VEP) task, we propose a multilevel multivariate functional principal components analysis (FPCA) for high-dimensional functional outcomes as a function of time and frequency. Applications to VEP data lead to new insights on autism-specific neural activity patterns. The autistic group shows significantly lower evoked and higher induced gamma power compared to the neurotypical group. In addition, while subject level variation is dominated by variation in the stimulus-locked evoked signal in neurotypical development, it is dominated by induced power in autism.

C1495: Canonical mutual information under meta-elliptic symmetry

Presenter: **Sarbojit Roy**, King Abdullah University of Science and Technology, Saudi Arabia

Understanding brain network connectivity requires interpretable models that can capture complex neural dependencies. While canonical correlation analysis is restricted to linear associations, more general dependence measures such as mutual information (MI) can detect both linear and nonlinear relationships but often lack interpretability at the covariate level. We propose a semiparametric, interpretable framework to quantify the global association between two brain regions using MI. Our network-based estimator is consistent under meta-elliptic symmetry of the covariates and can highlight key drivers of connectivity. These interpretability approaches are especially suited for neuroscience applications, where identifying and explaining connectivity patterns is essential.

C1496: Clustering spatial transcriptomics data with dirichlet process mixture of random spanning trees

Presenter: **Bani Mallick**, Texas A&M University, United States

Spatial transcriptomics has gained tremendous popularity as it allows researchers to map gene expression directly onto tissue architecture, preserving spatial context and providing high-resolution insights into cellular interactions and biological processes within their native environments. We introduce a novel Bayesian nonparametric framework, the Dirichlet process mixture of random spanning trees, designed to detect an unknown number of possibly non-convex clusters in possibly non-convex spatial domains. The model's two-layer partitioning effectively addresses challenges posed by the intricate spatial organization of tissue samples, such as non-convex clusters and irregular spatial boundaries of the samples. Simulation studies show that this method achieves superior clustering accuracy compared to existing methods. We apply our method to our motivating mouse colonic dataset during healing from inflammatory damage, revealing meaningful clusters associated with different stages of tissue repair.

CO087 Room BCB 311 EMERGING TOPICS IN STATISTICAL META-ANALYSIS AND APPLICATIONS

Chair: Din Chen

C0363: Geometric standardized mean difference and its application to meta-analysis

Presenter: **Tiejun Tong**, Hong Kong Baptist University, Hong Kong

The standardized mean difference (SMD) is a widely used measure of effect size, particularly common in psychology, clinical trials, and meta-analysis involving continuous outcomes. Traditionally, under the equal variance assumption, the SMD is defined as the mean difference divided by a common standard deviation. This approach is prevalent in meta-analysis but can be overly restrictive in clinical practice. To accommodate unequal variances, the conventional method averages the two variances arithmetically, which does not allow for an unbiased estimation of the SMD. Inspired by this, a geometric approach is proposed to average the variances, resulting in a novel measure for standardizing the mean difference with unequal variances. The Cohen-type and Hedges-type estimators are further proposed for the new SMD, and their statistical properties are derived, including the confidence intervals. Simulation results show that the Hedges-type estimator performs optimally across various scenarios, demonstrating lower bias, lower mean squared error, and improved coverage probability. A real-world meta-analysis also illustrates that the new SMD and its estimators provide valuable insights to the existing literature and can be highly recommended for practical use.

C0716: Statistical considerations for inverse publication bias in safety outcomes

Presenter: **Lifeng Lin**, University of Arizona, United States

Inverse publication bias (IPB) is an emerging concern in meta-analyses of safety outcomes, where studies favoring comparable safety profiles between interventions and controls may be more likely to be published. Unlike traditional publication bias, IPB is less recognized yet can lead to misleading conclusions in systematic reviews. Key statistical considerations are discussed for detecting IPB. First, methodological approaches are presented, including how contour-enhanced funnel plots and statistical tests such as Egger's and Peters' regressions can be adapted for IPB detection, emphasizing the impact of effect direction in bias assessments. Second, empirical analyses from meta-analyses of adverse events in the SMART Safety dataset, the largest evidence synthesis database for adverse events, are shared. Findings indicate that a considerable proportion of meta-analyses exhibit IPB and that qualitative assessments may be necessary to complement statistical methods, particularly in small-study contexts. The aim is to raise awareness of IPB and provide practical guidance for systematic reviews of safety outcomes.

C0755: Meta-analysis through combining p-values

Presenter: **Zhongxue Chen**, Arizona State University, United States

Many meta-analysis approaches, such as fixed- and random-effect models, make some distributional assumptions. If those assumptions are violated, a meta-analysis can still be conducted by combining p-values. Combining information (p-values) obtained from individual studies to test whether there is an overall effect is an important task in statistical data analysis. Many classical statistical tests, such as chi-squared tests, can be viewed as being a p-value combination approach. It remains challenging to find powerful methods to combine p-values obtained from various sources. In this project, a class of p-value combination methods is studied based on the gamma distribution. It is shown that this class of tests is optimal under certain conditions, and several existing popular methods are equivalent to its special cases. An asymptotically uniformly most powerful p-value combination test based on the constrained likelihood ratio test is then studied. Numeric results from the simulation study and real data examples demonstrate that the proposed tests are robust and powerful under many conditions. They have potential broad applications in statistical inference.

C0888: Confidence score: A data-driven measure for inclusive systematic reviews considering unpublished preprints

Presenter: **Jiayi Tong**, Johns Hopkins University, United States

Co-authors: Chongliang Luo, Yifei Sun, Rui Duan, M Elle Saine, Lifeng Lin, Yifan Peng, Yiwen Lu, Anchita Batra, Olivia Wang, Ruowang Li, Arielle Marks-Anglin, Yuchen Yang, Xu Zuo, Jiang Bian, Stephen Kimmel, Keith Hamilton, Adam Cuker, Rebecca Hubbard, Hua Xu, Yong Chen

COVID-19, since its emergence in December 2019, has globally impacted research, with over 360,000 COVID-19-related manuscripts published on PubMed and preprint servers like medRxiv and bioRxiv, where preprints comprise about 15% of all manuscripts; yet, the role and impact of preprints on COVID-19 research and evidence synthesis remain uncertain. A novel data-driven method is proposed for assigning weights to individual preprints in systematic reviews and meta-analyses, using a confidence score derived from the survival cure model (also known as the survival mixture model), which accounts for the time elapsed between posting and publication of a preprint, as well as metadata such as the number of first two-week citations, sample size, and study type. Using 146 preprints on COVID-19 therapeutics posted from the beginning of the pandemic through April 30, 2021, the confidence scores are validated, showing an area under the curve of 0.95 (95% CI, 0.92-0.98), and through a use case on the effectiveness of hydroxychloroquine, demonstrating how these scores can be practically incorporated into meta-analyses to properly weigh preprints.

C0870: From big-data to small area estimation: Statistical meta-analysis approach

Presenter: **Din Chen**, University of Pretoria, South Africa

Co-authors: Jeffrey Wilson

In the era of big data deluge, meta-analytics offers unprecedented opportunities to enhance the precision and relevance of statistical inference, especially in public health and social sciences. Another critical area of meta-analytics is small area estimation (SAE), where reliable statistics are required at granular geographic or demographic levels, often with limited local data. The purpose is to explore the integration of big data analytics and SAE with meta-analytic frameworks. It begins by discussing the challenges in big-data sciences and SAE, and then by introducing a meta-analysis-based approach as a flexible tool for synthesizing diverse data sources. The methodological innovations, computational strategies, and practical implications for policy and planning are highlighted. This cross-disciplinary approach provides a blueprint for leveraging big data to inform local decision-making through statistically principled estimation.

CO051 Room BCB 312 DESIGN AND ANALYSIS OF EXPERIMENTS

Chair: Steven Gilmour

C1002: Optimal designs for a linear paired comparison model with continuous predictors in a restricted design region

Presenter: **Heiko Grossmann**, Otto-von-Guericke-University Magdeburg, Germany

Co-authors: Heinz Holling, Nadja Malevich, Rainer Schwabe

A linear model is considered with a continuous response for pairwise comparisons of alternatives characterized by several continuous factors. For each pair of alternatives, the response may then be interpreted as the difference between the expected utilities of the two alternatives, plus a random error. Other interpretations arise from the fact that the model is formally equivalent to a multiple linear regression model without an intercept, where the factor levels of one alternative in each pair enter the systematic component of the model additively, while the factor levels of the other are included in a subtractive manner. For example, this corresponds to experimental settings in which two simultaneously applied factors exert opposing effects on the outcome for each experimental unit. In standard design problems, the design region consists of all pairs of vectors in the Cartesian product of certain intervals. However, a restricted design region and an associated design locus are considered, both of which are the union of certain lower-dimensional sets. Optimal continuous designs and efficiency comparisons are presented for various criteria, including D-, A-, and MV-optimality. Illustrative examples of potential applications of the model and the proposed designs are outlined.

C0499: An adaptive regrouping subsampling method for linear mixed models

Presenter: **Rundong Zhou**, King's College London, United Kingdom

In latent group modeling scenarios, datasets are often assumed to follow a linear mixed model, although group memberships are unknown. In addition, group information may be missing or unreliable, which has been largely overlooked in previous studies. A key assumption is introduced: Grouping information provided in a dataset, particularly in linearly structured data, should not be fully trusted. Instead, mathematical procedures are applied to determine whether the data conform to a single linear model or a linear mixed model, and to identify appropriate groupings within a linear mixture framework. To support this, an unsupervised subsampling algorithm is developed that incorporates clustering to automatically assign group labels without requiring prior grouping information. This method ensures that the model fitting process retains the most informative results under potential group uncertainty.

C1108: Maximum likelihood estimation under the Emax model

Presenter: **Caterina May**, Università del Piemonte Orientale, Italy

Co-authors: Chiara Tommasi, Giacomo Aletti, Nancy Flournoy

The purpose is to examine the estimation of the Emax dose response model, a commonly employed framework across various fields including clinical trials, pharmacology, agriculture, and environmental studies. Although difficulties in computing maximum likelihood estimates (MLEs) for the model parameters of the Emax model are frequently attributed to numerical or computational limitations, they more accurately arise from the non-existence of an MLE in certain cases. The contribution offers new insights and practical guidance for handling the full range of experimental scenarios that practitioners may encounter, supporting them throughout the estimation process.

C1196: Orthogonal factorial designs for trials of therapist-delivered interventions*Presenter:* **Rebecca Walwyn**, University of Leeds, United Kingdom*Co-authors:* Rosemary Bailey, Arpan Singh, Neil Corrigan, Steven Gilmour

It is recognised that treatment-related clustering should be allowed for in the sample size and analyses of individually-randomised parallel-group trials evaluating therapist-delivered interventions such as psychotherapy. Interventions are a treatment factor, but therapists are not. If the aim of a trial is to separate effects of therapists from those of interventions, interventions and therapists are regarded as two potentially interacting treatment factors (one fixed, one random) with a factorial structure. The specific design is considered where each therapist delivers each intervention (crossed therapist-intervention design), and the resulting therapist-intervention combinations are randomised to patients. A classical design of experiments approach is adopted to propose a family of orthogonal factorial designs and their associated data analyses, which also allow for therapist learning and centre. The associated data analyses are set out using ANOVA and regression and report the results of a small simulation study that was conducted to explore the performance of the proposed randomization methods on estimating the intervention effect, the between-therapist variance, and the between-therapist variance in the intervention effect. It is concluded that a more purposeful trial design has the potential to lead to better evidence on a range of complex healthcare interventions and outline areas for further methodological research.

| | |
|--|---------------------------------|
| CO209 Room BCB 402 BAYESIAN METHODS IN ENVIRONMENTAL, EPIDEMIOLOGIC, GENOMIC PROBLEMS | Chair: Andrea Sottosanti |
|--|---------------------------------|

C0605: Leveraging satellite imagery to estimate number of households at fine scale resolutions in resource-poor settings*Presenter:* **Chibuzor Christopher Nnanatu**, University of Southampton, United Kingdom

Effective government policies and humanitarian response efforts require accurate knowledge of the population count and the number of households at small area scales. In settings where census data are unavailable, incomplete, or outdated, advanced statistical models have been developed to produce estimates of counts of people at fine spatial resolutions, but methods for producing estimates of number of households are generally lacking. A statistical modeling technique is presented for producing fine spatial resolution estimates of the number of households by integrating AI-powered satellite-derived human settlement maps with stacks of geospatial datasets and demographic datasets. Parameter estimation is based on a Bayesian statistical inference approach implemented via the integrated nested Laplace approximation with stochastic partial differential equation frameworks, thereby making uncertainty quantification straightforward. The methodology is evaluated using a simulation study, showing varying levels of accuracy over different magnitudes of data missingness, with lower estimation error obtained for a smaller proportion of missing samples. The approach was successfully applied to obtain estimates of the number of households along with the corresponding estimates of uncertainty at 100m grid cells across Cameroon. The methodology provides a significant advancement in the small area population estimation field for facilitating more efficient governance and resource allocation.

C0625: Graphical models for modern biological applications*Presenter:* **Francesco Stingo**, University of Florence, Italy

The purpose is to discuss recent inferential and computational techniques for multiple graphical models, where the sub-group assignment depends on the value of an externally observed covariate. Bayesian Gaussian graphical models with covariates (GGMx) are then introduced, a class of multivariate Gaussian distributions with covariate-dependent sparse precision matrix. A general construction of a functional mapping is proposed from the covariate space to the cone of sparse positive definite matrices encompassing many existing graphical models for heterogeneous settings. The flexible formulation of GGMx allows both the strength and the sparsity pattern of the precision matrix (hence the graph structure) to change with the covariates. Extensive simulations and a cancer genomics case study demonstrate the proposed models' utility.

C0858: Leveraging drug therapeutic class to identify adverse events from drug-drug interactions: A decision-theoretic approach*Presenter:* **Massimiliano Russo**, the Ohio State University, United States

The concurrent prescription of multiple drugs can increase the likelihood of adverse drug events due to potential interactions. Identifying these adverse events is complex, especially when considering rare events and infrequently prescribed drug combinations. A novel Bayesian framework is introduced, designed to detect potential adverse events resulting from drug-drug interactions by incorporating information about each drug's therapeutic class (THERCL) to enhance signal detection. Using a decision-theoretic framework facilitates the identification of potentially harmful drug combinations while explicitly controlling the proportion of false-positive results. It is shown that the proposed framework outperforms the state-of-the-art approaches in an extensive simulation study and application involving a cohort study of older adults. Results indicate that leveraging THERCL in Bayesian modeling can effectively reduce false positives in detecting adverse events from drug-drug interactions, potentially improving safety in medication prescribing.

C1060: A Bayesian spatiotemporal Poisson auto-regressive model for the disease infection rate*Presenter:* **Pierfrancesco Alaimo Di Loro**, LUMSA University, Italy*Co-authors:* Sujit Sahu, Dankmar Boehning

The COVID-19 pandemic provided new modelling challenges to investigate epidemic processes. Poisson auto-regression is extended to incorporate spatiotemporal dependence and characterize the local dynamics by borrowing information from adjacent areas. Adopted in a fully Bayesian framework and implemented through a novel sparse-matrix representation in Stan, the model has been validated through a simulation study. It is used to analyze the weekly COVID-19 cases in the English local authority districts and verify some of the epidemic-driving factors. The model detects substantial spatiotemporal heterogeneity and enables the formalization of novel model-based investigation methods for assessing additional aspects of disease epidemiology.

| | |
|--|-------------------------|
| CO225 Room BCB 403 ADVANCES IN STATISTICAL ANALYSES OF NETWORK DATA | Chair: Yinqiu He |
|--|-------------------------|

C0541: SANVI: A fast spectral-assisted network variational inference method with an extended surrogate likelihood function*Presenter:* **Fangzheng Xie**, Indiana University, United States*Co-authors:* Dingbo Wu, Fangzheng Xie

Bayesian inference has been broadly applied to statistical network analysis, but suffers from the expensive computational costs due to the nature of Markov chain Monte Carlo sampling algorithms. A novel and computationally efficient spectral-assisted network variational inference (SANVI) method is proposed under the generalized random dot product graph framework. The key idea is a cleverly designed extended surrogate likelihood function that enjoys two convenient features. Firstly, it decouples the generalized inner product of latent positions in the random graph model. Secondly, it extends the domain of the original likelihood function to the entire Euclidean space. Leveraging these features, a computationally efficient Gaussian variational inference algorithm is designed via the stochastic gradient descent method for Bayesian inference of networks. Furthermore, the asymptotic efficiency of the maximum extended surrogate likelihood estimator and the Bernstein-von Mises limit of the variational posterior distribution is shown. Through extensive numerical studies, the practical advantage of the proposed SANVI algorithm compared to the classical Markov chain Monte Carlo algorithm is demonstrated, including the estimation accuracy for the latent positions and the computational costs.

C0717: Adaptive robust confidence intervals: Location models and networks*Presenter:* **Yuetian Luo**, Rutgers University, United States*Co-authors:* Chao Gao

The purpose is to study the construction of confidence intervals under Huber's contamination model. When the contamination proportion is

unknown, the necessary adaptation cost of the problem is characterized. In particular, for the Gaussian location model, the optimal length of an adaptive confidence interval is proved to be exponentially wider than that of a non-adaptive one. Results for general location models will be discussed. In addition, the same problem is also considered in a network setting for an Erdos-Renyi graph with node contamination. It is shown that the hardness of the adaptive confidence interval construction is implied by the detection threshold between the Erdos-Renyi model and the stochastic block model.

C0906: Conformal prediction for dyadic regression

Presenter: **Robert Lunde**, Washington University in St Louis, United States

Co-authors: Ji Zhu, Liza Levina

Dyadic regression, which involves modeling a relational matrix given covariate information, is an important task in statistical network analysis. Uncertainty quantification is considered for dyadic regression models using conformal prediction. Finite-sample validity of the procedures is established for various sampling mechanisms under a joint exchangeability assumption. The proof uses new results related to the validity of conformal prediction beyond exchangeability, which may be of independent interest. It is also shown that, under certain conditions, it is possible to construct asymptotically valid prediction sets for a missing entry under a structured missingness assumption.

C1119: Principal graph encoder embedding and principal community detection

Presenter: **Yuexiao Dong**, Temple University, United States

Co-authors: Cencheng Shen, Carey Priebe, Youngser Park

The purpose is to introduce the concept of principal communities and propose a principal graph encoder embedding method that concurrently detects these communities and achieves vertex embedding. Given a graph adjacency matrix with vertex labels, the method computes a sample community score for each community, ranking them to measure community importance and estimate a set of principal communities. The method then produces a vertex embedding by retaining only the dimensions corresponding to these principal communities. Theoretically, we define the population version of the encoder embedding and the community score based on a random Bernoulli graph distribution. It is proven that the population principal graph encoder embedding preserves the conditional density of the vertex labels and that the population community score successfully distinguishes the principal communities. A variety of simulations is conducted to demonstrate the finite-sample accuracy in detecting ground-truth principal communities, as well as the advantages in embedding visualization and subsequent vertex classification. The method is further applied to a set of real-world graphs, showcasing its numerical advantages, including robustness to label noise and computational scalability.

CO070 Room BCB 405 MONTE CARLO METHODS AND THEIR APPLICATION

Chair: Galin Jones

C0197: SOMA: A novel sampler for exchangeable variables

Presenter: **Nianqiao Ju**, Dartmouth College, United States

Co-authors: Yifei Xiong

The problem of sampling exchangeable random variables arises in many Bayesian inference tasks, especially in data imputation given a privatized summary statistic. These permutation-invariant joint distributions often have dependency structures that make sampling challenging. Component-wise sampling strategies, such as Metropolis-within-Gibbs, can mix slowly because they consider only comparing a proposed point with one component at a time. A novel single-offer-multiple-attempts (SOMA) sampler is proposed that is tailored to sampling permutation-invariant distributions. The core intuition of SOMA is that a proposed point unsuitable to replace one component might still be a good candidate to replace some other component in the joint distribution. SOMA first makes a single offer, and then simultaneously considers attempts to replace each component of the current state with the single offer, before making the final acceptance or rejection decision. An acceptance lower bound of SOMA is provided, and using a coupling method, the convergence rate upper bound of SOMA is derived in the two-dimensional case. Theoretical findings are validated with numerical simulations, including a demonstration of differentially private Bayesian linear regression.

C0592: Some insights into the reliability and limitations of adaptive MCMC

Presenter: **Austin Brown**, Texas AM University, United States

The reliability and limitations of adaptation strategies used in adaptive Markov chain Monte Carlo (MCMC) are investigated. In particular, general lower bounds are established on the weak convergence rate under general adaptation plans. If the adaptation diminishes sufficiently fast, comparable convergence rate upper bounds are also developed. These results provide some insight into the optimal design of adaptation strategies and help set expectations for the practical performance of adaptive MCMC. Applications to an adaptive unadjusted Langevin algorithm and to adaptively tuning the covariance matrix in random-walk Metropolis-Hastings are explored.

C1055: Statistical properties of initial sequence type variance estimators for reversible Markov chains

Presenter: **Hyebin Song**, The Pennsylvania State University, United States

Initial sequence estimators, originally introduced by a prior study, are commonly used to estimate Monte Carlo standard errors in reversible Markov chain Monte Carlo chains. In particular, the initial positive sequence estimator utilizes the property that the sums of adjacent autocovariances are non-negative, summing autocovariances up to the point where this non-negativity condition is violated. While this estimator has been widely used and adapted to different settings, only its asymptotic conservativeness has been shown, while consistency remains an open question. This gap is addressed by investigating the statistical properties of the initial positive sequence estimator. The convergence behavior of its random truncation point is first studied. An alternative initial sequence-type estimator is also proposed based on a modified truncation rule. For both estimators, consistency is established, and bounds are derived on rates of convergence. Finally, through empirical studies using both simulated and real-world data, the theoretical findings are validated, and the empirical performance of the two initial sequence-type estimators is compared with the standard overlapping batch mean estimator.

C1441: Adaptive stereographic MCMC

Presenter: **Krzysztof Latuszynski**, University of Warwick, United Kingdom

Co-authors: Cameron Bell, Gareth Roberts

In order to tackle the problem of sampling from heavy-tailed, high-dimensional distributions via Markov chain Monte Carlo (MCMC) methods, an earlier work introduces Stereographic MCMC samplers. However, its improvement in algorithmic efficiency, as well as the computational cost of the implementation, is significantly impacted by the parameters used in this design. To address these design difficulties, the adaptive versions of three stereographic MCMC algorithms - the stereographic random walk (SRW), the stereographic slice sampler (SSS), and the stereographic bouncy particle sampler (SBPS) - is introduced, which automatically update the parameters of the algorithms as the run progresses. The adaptive setup allows for better exploitation of the power of the stereographic projection, even when the target distribution is neither centered nor homogeneous. Unlike Hamiltonian Monte Carlo (HMC) and other off-the-shelf MCMC samplers, the resulting algorithms are robust to starting far from the mean in heavy-tailed, high-dimensional settings. To prove convergence properties, a novel framework is developed for the analysis of adaptive MCMC algorithms over collections of simultaneously uniformly ergodic Markov operators, which is applicable to continuous-time processes, such as SBPS. This framework allows obtaining L2 and almost sure convergence results, and a CLT for the adaptive stereographic algorithms.

CO103 Room BCB 406 NETWORK ANALYSIS METHODS AND THEIR APPLICATIONS**Chair: Panpan Zhang****C0479: Generalized Bayesian inference for dynamic random dot product graphs****Presenter:** Joshua Loyal, Florida State University, United States

The random dot product graph is a popular model for network data with extensions that accommodate dynamic (time-varying) networks. However, two significant deficiencies exist in the dynamic random dot product graph literature: (1) no coherent Bayesian way to update one's prior beliefs about the model parameters due to their complicated constraints, and (2) no approach to forecast future networks with meaningful uncertainty quantification. A generalized Bayesian framework is proposed that addresses these needs using a Gibbs posterior that represents a coherent updating of Bayesian beliefs based on a least-squares loss function. The consistency and contraction rate of this Gibbs posterior are established under commonly adopted Gaussian random walk priors. For estimation, a fast Gibbs sampler is developed with a time complexity that is linear in both the number of time points and observed edges in the dynamic network. Simulations and real data analyses show that the proposed methods in-sample and forecasting performance outperform that of competitors.

C0523: Accounting for network noise in graph-guided Bayesian modeling of -omics data**Presenter:** Wenrui Li, University of Connecticut, United States**Co-authors:** Changge Chang, Suprateek Kundu, Qi Long

There is a growing body of literature on knowledge-guided statistical learning methods for analysis of -omics data that can incorporate knowledge of underlying networks derived from functional genomics and functional proteomics. These methods have been shown to improve variable selection and prediction accuracy, and yield more interpretable results. However, these methods typically use graphs extracted from existing databases or rely on subject matter expertise, which are known to be incomplete and may contain false edges. To address this gap, a graph-guided Bayesian modeling framework is proposed to account for network noise in regression models involving structured high-dimensional predictors. Specifically, two sources of network information are used, including the noisy graph extracted from existing databases and the estimated graph from observed predictors in the dataset at hand, to inform the model for the true underlying network via a latent scale modeling framework. This model is coupled with the Bayesian regression model with structured high-dimensional predictors involving an adaptive structured shrinkage prior. An efficient Markov chain Monte Carlo algorithm is developed for posterior sampling. The advantages of the method are demonstrated over existing methods in simulations, and through analyses of a genomics dataset and another proteomics dataset for Alzheimer's disease.

C0552: A joint modeling approach for multilayer egocentric social networks, with application to mental health studies**Presenter:** Kehui Chen, University of Pittsburgh, United States**Co-authors:** Xiaolin Bo

A joint modeling framework is proposed that simultaneously models the ego's response and their multi-layer, longitudinal egocentric networks. The model consists of a set of mixed response models for the ego's response and ego-alter ties linked through shared random effects. It is discussed how to account for the heterogeneity in egocentric networks and how to address the computational challenges posed by the mixed responses and multivariate random effects. Data analysis reveals that perceived closeness with social ties is significantly associated with depressive symptoms.

C0732: Genetic regression analysis of structure-function connectome coupling**Presenter:** Eardi Lila, University of Washington, United States

Magnetic resonance imaging has significantly improved the understanding of the connectivity patterns within the human brain by enabling measurement of the strength of anatomical connections between brain regions through white matter fibers (structural connectivity) and the degree of coactivation of brain regions (functional connectivity). Heritability analyses of connectivity have established that genetics account for a considerable portion of the observed intersubject variability. However, such analyses typically ignore the multidimensional nature of functional and structural connectomes. Observed brain connectivity is modeled as the sum of multidimensional latent genetic and environmental contributions, and a novel constrained estimator is introduced for the covariance matrices of the genetic and environmental components. The estimator is several orders of magnitude faster than existing methods without sacrificing estimation accuracy. The proposed covariance estimate provides a summary statistic that can be used to estimate the parameters of a novel regression analysis that enables the characterization of the relationship between the latent genetic components of structural and functional connectomes. The analysis suggests that the genetic component of functional connectomes is more highly predictable from the genetic component of structural connectomes, suggesting a close relationship at the genetic level that is attenuated by distinct environmental factors.

CO285 Room BCB 407 INNOVATIONS IN CAUSAL, PREDICTIVE, AND NONPARAMETRIC INFERENCE**Chair: Bingkai Wang****C0237: How to achieve model-robust inference in stepped wedge trials with model-based methods?****Presenter:** Bingkai Wang, University of Michigan, United States

A stepped wedge design is an unidirectional crossover design where clusters are randomized to distinct treatment sequences. While model-based analysis of stepped wedge designs is a standard practice to evaluate treatment effects accounting for clustering and adjusting for covariates, their properties under misspecification have not been systematically explored. The focus is on model-based methods, including linear mixed models and generalized estimating equations with an independence, simple exchangeable, or nested exchangeable working correlation structure. The purpose is to study when a potentially misspecified working model can offer consistent estimation of the marginal treatment effect estimands, which are defined nonparametrically with potential outcomes and may be functions of calendar time and/or exposure time. A central result is proven, that consistency for nonparametric estimands usually requires a correctly specified treatment effect structure, but generally not the remaining aspects of the working model (functional form of covariates, random effects, and error distribution), and valid inference is obtained via the sandwich variance estimator. Furthermore, an additional g-computation step is required to achieve model-robust inference under non-identity link functions or for ratio estimands. The theoretical results are illustrated via several simulation experiments and re-analysis of a completed stepped wedge cluster randomized trial.

C0252: On robust empirical likelihood for nonparametric regression with application to regression discontinuity designs**Presenter:** Qin Fang, the University of Sydney, Australia**Co-authors:** Shaojun Guo, Xinghao Qiao

Empirical likelihood serves as a powerful tool for constructing confidence intervals in nonparametric regression and regression discontinuity designs (RDD). The original empirical likelihood framework can be naturally extended to these settings using local linear smoothers, with Wilks' theorem holding only when an undersmoothed bandwidth is selected. However, the generalization of bias-corrected versions of empirical likelihood under more realistic conditions is non-trivial and has remained an open challenge in the literature. A satisfactory solution is provided by proposing a novel approach, referred to as robust empirical likelihood, designed for nonparametric regression and RDD. The core idea is to construct robust weights that simultaneously achieve bias correction and account for the additional variability introduced by the estimated bias, thereby enabling valid confidence interval construction without extra estimation steps involved. It is demonstrated that the Wilks' phenomenon still holds under weaker conditions in nonparametric regression, sharp and fuzzy RDD settings. Moreover, the proposed procedure exhibits robustness to bandwidth selection, making it a flexible and reliable tool for empirical analyses. The practical usefulness is illustrated through extensive simulations and applications to two real datasets.

C0608: SymmPI: Predictive inference for data with group symmetries*Presenter:* **Mengxin Yu**, Washington University in St. Louis, United States

Quantifying the uncertainty of predictions is a core problem in modern statistics. Methods for predictive inference have been developed under a variety of assumptions, often, for instance, in standard conformal prediction, relying on the invariance of the distribution of the data under special groups of transformations such as permutation groups. Moreover, many existing methods for predictive inference aim to predict unobserved outcomes in sequences of feature-outcome observations. Meanwhile, there is interest in predictive inference under more general observation models (e.g., for partially observed features) and for data satisfying more general distributional symmetries (e.g., network, rotationally invariant, or coordinate-independent observations in physics). SymmPI is proposed, a unified methodology for predictive inference when data distributions have general group symmetries in arbitrary observation models. The methods leverage the novel notion of distributional equivariant transformations, which process the data while preserving their distributional invariances. It is shown that SymmPI has valid coverage under distributional invariance, and its performance is characterized under distribution shift, recovering recent results as special cases. These methodologies are particularly relevant for cluster-randomized trials in clinical settings, where prediction reliability is essential.

C0613: Proximal causal inference for conditional separable effects*Presenter:* **Chan Park**, University of Illinois, Urbana-Champaign, United States*Co-authors:* Mats Stensrud, Eric Tchetgen Tchetgen

Scientists regularly pose questions about treatment effects on outcomes conditional on a post-treatment event. However, defining, identifying, and estimating causal effects conditional on post-treatment events requires care, even in perfectly executed randomized experiments. Recently, the conditional separable effect (CSE) was proposed as an interventionist estimand that corresponds to scientifically meaningful questions in these settings. However, while being a single-world estimand, which can be queried experimentally, existing identification results for the CSE require no unmeasured confounding between the outcome and post-treatment event. This assumption can be violated in many applications. This concern is addressed by developing new identification and estimation results for the CSE in the presence of unmeasured confounding. Nonparametric identification of the CSE is established in observational and experimental settings when time-varying confounders are present, and certain proxy variables are available for hidden common causes of the post-treatment event and outcome. For inference, an influence function is characterized for the CSE under a semiparametric model in which nuisance functions are a priori unrestricted. Moreover, a consistent, asymptotically linear, and locally semiparametric efficient estimator of the CSE is developed using modern machine learning theory. The framework is illustrated with simulation studies and a real-world cancer therapy trial.

CO118 Room BCB 408 BAYESIAN SEMI- AND NON-PARAMETRIC METHODS I**Chair: Guillaume Kon Kam King****C0463: Discrete autoregressive switching processes in sparse graphical modeling of multivariate time series data***Presenter:* **Beniamino Hadj-Amar**, University of South Carolina, United States*Co-authors:* Aaron Bornstein, Michele Guindani, Marina Vannucci

A flexible Bayesian approach is proposed for sparse Gaussian graphical modeling of multivariate time series. Temporal correlation is accounted for in the data by assuming that observations are characterized by an underlying and unobserved hidden discrete autoregressive process. Multivariate Gaussian emission distributions are assumed, and spatial dependencies are captured by modeling the state-specific precision matrices via graphical horseshoe priors. The mixing probabilities of the hidden process are characterized via a cumulative shrinkage prior that accommodates zero-inflated parameters for non-active components, and further incorporates a sparsity-inducing Dirichlet prior to estimate the effective number of states from the data. For posterior inference, a sampling procedure is developed that allows estimation of the number of discrete autoregressive lags and the number of states, and that cleverly avoids having to deal with the changing dimensions of the parameter space. Performance of the proposed methodology is thoroughly investigated through several simulation studies. The use of the approach is further illustrated for the estimation of dynamic brain connectivity based on fMRI data collected on a subject performing a task-based experiment on latent learning

C0706: Latent position co-clustering for multiple social networks*Presenter:* **Michael Fop**, University College Dublin, Ireland*Co-authors:* CJ Clarke

Social networks often involve multiple types of relationships, such as friendship, collaboration, and communication, resulting in a collection of networks recording different social dimensions over the same set of individuals. Analyzing such data requires methods that can uncover structure within single networks and across different relational dimensions. A latent position co-clustering model is introduced that jointly clusters networks and their constituent nodes. Built on a hierarchical mixture-of-mixtures formulation, the model simultaneously performs dimension reduction and two-level clustering. At the network level, it groups networks that share similar latent topological structures. At the node level, it uncovers local connectivity patterns such as communities or social roles. The latent space provides a parsimonious and interpretable representation of both global and local structure. The model adopts a Bayesian nonparametric framework based on mixtures of finite mixtures, which place priors on the number of mixture components at both levels and incorporate sparse priors to encourage parsimonious clustering. Inference is conducted via Markov chain Monte Carlo, employing a telescoping sampling strategy and a tailored post-processing procedure. Applications to real-world social multiplex data reveal interpretable network-level clusters aligned with contextual features, and node-level clusters that reflect roles and social patterns.

C0966: Bayesian nonparametric models with BART components*Presenter:* **Maria Kalli**, Kings College London, United Kingdom*Co-authors:* Jim Griffin

Bayesian additive regression tree, BART, models have emerged as an important method for Bayesian nonlinear regression. However, there is little work on fully nonparametric versions of BART. Bayesian nonparametric models are described, built using BART elements, which can be easily implemented using existing BART R packages. A general density regression method is developed using normalized latent measure random factor models to build mixtures where the components are heteroscedastic BARTs. The performance of the method is evaluated on both simulated and real data.

C0879: Calibrating covariates within product partition models, with an application to stochastic block models*Presenter:* **Sirio Legramanti**, University of Bergamo, Italy*Co-authors:* Raffaele Argiento, Valentina Ghidini

Product partition models (PPM) represent a flexible framework for Bayesian nonparametric clustering. Their factorized structure facilitates the incorporation of individual covariates, giving rise to variants like the well-known PPM with covariates (PPMx) and spatial PPM (sPPM), with the latter being specialized to spatial covariates. Besides incorporating covariates into the clustering process, it is paramount to calibrate their influence on the obtained partition. A framework is proposed for weighting the influence of covariates within PPMs, and findings are illustrated on stochastic block models. The latter are models for clustering network nodes based on the network adjacency matrix. Such an application is further motivated by the fact that network data are often accompanied by node covariates. For example, in the real-data application to the public transportation network of Bergamo province (Italy), each network node has a spatial location, and one may aim for clusters that are as spatially cohesive as possible, while still reflecting the network structure. In fact, the obtained clusters may be used to inform policymaking in public transport, and it may be preferable that such policies are uniform over neighboring areas.

| | |
|---|------------------------------|
| CO091 Room BCB 409 STATISTICAL METHODS AND APPLICATION FOR HIGH DIMENSIONAL BIOMARKERS | Chair: Dana Tudorascu |
|---|------------------------------|

C0750: Biclustering multivariate longitudinal data in sport-related concussion recovery*Presenter:* **Jaroslav Harezlak**, Indiana University School of Public Health-Bloomington, United States*Co-authors:* Luo Xiao, Yu-Chien Wu, Qiuting Wen, Caleb Weaver

Biclustering involves the simultaneous clustering of both samples and features within a dataset, enabling the identification of subsets of samples that exhibit similar behavior across specific subsets of features. Motivated by a longitudinal diffusion tensor imaging (DTI) study of sport-related concussion (SRC), the problem of biclustering multivariate longitudinal data is addressed. In this context, subjects and features are grouped based on shared longitudinal patterns rather than absolute magnitudes. A penalized regression-based method is proposed that leverages heterogeneity in these temporal patterns across subjects and features. The effectiveness of the approach is demonstrated through a simulation study and an application to the motivating DTI dataset. The analysis uncovers distinct subgroups of SRC cases characterized by heterogeneous patterns of white matter abnormalities.

C0752: Multimodal topological analysis of neuroimaging in Alzheimer's disease via persistent homology*Presenter:* **Yuexuan Wu**, University of South Carolina, United States

Alzheimer's disease (AD) involves complex pathological processes, including the accumulation of beta-amyloid and tau proteins and progressive structural brain degeneration. Early and accurate prediction of disease onset and progression remains a critical challenge. A unified persistent homology framework is proposed to analyze multimodal neuroimaging data for both mechanistic understanding and predictive modeling. A bi-filtered approach is first introduced for positron emission tomography (PET), enabling threshold-free characterization of the spatial co-localization and interaction between beta-amyloid and tau. By embedding PET-derived brain networks into a dendrogram-based geometric space, multiscale topological patterns that reveal biologically meaningful features are captured, and thus, a stronger association with clinical AD. Additionally, topological features of brain morphology are extracted from structural MRI via cubical persistent homology, and these features are used to reveal patterns of neurodegeneration, classify disease stages, and forecast progression. The integration of topological features across multiple imaging modalities improves disease prediction performance over existing methods. Together, the utility of the persistent homology-based approach is demonstrated in neuroimaging for both biomarker discovery and robust prediction, offering valuable insights for the early detection of AD.

C0772: A comparative study of digital biomarkers of variability in glucose levels among individuals with type I diabetes*Presenter:* **Adi Andrei**, Northwestern University, United States

Type I diabetes (T1D) is an autoimmune condition characterized by the destruction of insulin-producing beta cells in the pancreas, leading to a lifelong dependency on insulin therapy. Continuous glucose monitoring (CGM) technology has emerged as a valuable tool for evaluating and maintaining adequate glycemic control in individuals with T1D. While maintaining good average glycemic control is important, it is not sufficient to robustly characterize the risk profile of individuals with T1D. Despite similar average blood glucose levels, individuals may present vastly different glucose variability patterns, leading to different clinical risks and outcomes. Day-to-day variability metrics inform how much glucose patterns differ from one day to the next and, unlike variability measures that focus on fluctuations within a single day, assess the reproducibility of glucose control across multiple days. Consequently, such metrics reveal whether a patient's diabetes management is consistently effective and may inform treatment optimization and insulin dose adjustments. Leveraging CGM data collected in a recent T1D clinical trial, daily glucose variability measures are quantified and compared based on cumulative distribution functions. Findings suggest that strategies aimed at reducing the day-to-day glucose level variability have the potential to result in superior personalized glycemic control clinical recommendations.

C0782: Federated factor analysis for collaborative learning of multi-site imaging data*Presenter:* **Qiong Wu**, University of Pittsburgh, United States*Co-authors:* Yezhi Pan, Shuo Chen

Multi-site collaboration is increasingly common in biomedical science, providing valuable opportunities to analyze complex datasets, such as neuroimaging data, that require extensive time and resources for collection. Factor analysis is a widely used statistical technique to uncover multivariate relationships by identifying latent factors. However, re-identification risks and privacy policies restrict the sharing of sensitive individual-level data across sites, challenging multi-site factor analysis. Traditional per-site factor analysis can yield inconsistent definitions of latent factors, complicating interpretability when aggregating results across sites. Another primary challenge is the inherent data heterogeneity arising from variations in patient populations and technical factors across sites. Hence, a structure-guided confirmatory factor analysis (SCFA) is proposed to identify a unified set of latent factors across sites while allowing site-specific heterogeneity in factor covariances. A federated learning algorithm, Fed-Factor, is introduced to solve the SCFA model securely with only a single round of summary statistic communication across sites, producing results identical to those obtained from pooling individual-level data. Its effectiveness and utility are demonstrated through simulations and imaging data from Adolescent Brain Cognitive Development (ABCD) study.

| |
|--|
| CC442 Room BCB 213 SHORT TALKS: CMSTATISTICS II |
|--|

| |
|------------------------------------|
| Chair: Nilanjan Chakraborty |
|------------------------------------|

C1431: Estimating semiparametric Gaussian mixtures of nonparametric regressions with an application in environmental economics*Presenter:* **Sphive Skhosana**, University of Pretoria, South Africa*Co-authors:* Sollie Millard, Frans Kanfer

Semiparametric mixtures of Gaussian nonparametric regressions (SPGMNRs) are a flexible class of Gaussian mixtures of regression models. These models assume that the component regression functions (CRFs) are nonparametric functions of the covariates, whereas the mixing proportions and variances are constant (parametric). However, local-likelihood estimation of the nonparametric (CRFs) poses a computational challenge. Traditional expectation-maximization (EM) optimisation of the local-likelihood functions is not appropriate due to the label-switching problem. Separately applying the EM algorithm on each local-likelihood function will likely result in non-smooth CRFs. This is because the local responsibilities calculated at the local E-step of each EM are not guaranteed to be aligned. The misalignment in the labels can be prevented by making use of the same (global) responsibilities at each local M-step. Thus, the goal is to obtain these global responsibilities. The aim is to propose a novel two-step approach to address label-switching. In the first step, each local-likelihood function is maximized separately to obtain the local responsibilities. In the second step, based on an appropriate objective function, one set of local responsibilities is chosen from the first step as the global responsibilities. The performance and practical usefulness of the proposed method are evaluated using a simulated dataset and a real dataset from environmental economics.

C1504: Sample-specific Multiomic Association networks Using Gaussian graphical models*Presenter:* **Enakshi Saha**, University of South Carolina, United States

Gaussian Graphical Models (GGM) provide an invaluable tool for studying the interaction network between multiple omics modalities. However, existing methods estimate a single network that approximates the average conditional dependence structure across the entire population and fail to recognize individual-specific variability. To overcome this limitation, we propose an empirical Bayesian model, SMAUG (Sample-specific Multiomic Association networks Using Gaussian graphical models), that recognizes individual-specific heterogeneity in molecular dependence by estimating sample-specific GGMs. By employing data-driven Individual-specific conjugate priors, SMAUG provides a scalable tool for deciphering variability in disease mechanisms across sex, age and other clinical variables, thereby providing a more nuanced understanding of diseases. In addition, SMAUG, being a partial-correlation-based method, is better suited to distinguish between direct molecular dependence and spurious

correlations, compared to existing methods for sample-specific network inference that employ Pearson's correlation as their foundation. We demonstrate the efficacy of SMAUG using simulated and real datasets.

C1510: **Second-order sparse sufficient dimension reduction with applications to quadratic discriminant Analysis**

Presenter: **Jing Zeng**, University of Science and Technology of China, China

Motivated by exploratory data analysis, sufficient dimension reduction (SDR) methods, especially inverse regression methods such as sliced inverse regression (SIR) and sliced averaged variance estimation (SAVE), have been central to multivariate analysis for more than three decades. Despite their popularity, extending these methods to high-dimensional settings remains challenging. This paper addresses the computational and theoretical limitations of the less explored second-order SDR methods in high dimensions. We introduce a novel approach for sparse subspace estimation that utilizes quadratic convex optimization and leverages the group structure of tensor parameters, achieving significant parameter reduction. The proposed two-step estimator achieves consistency in dimension selection, variable selection, and subspace estimation at a high convergence rate under mild conditions. The effectiveness and efficiency of the proposed method are further demonstrated through extensive simulation studies and real data examples. Additionally, the proposed sparse second-order SDR techniques are applied to quadratic discriminant analysis (QDA) problems and provide practitioners with a sparse projective classification method that is theoretically guaranteed and empirically well-performed.

C1512: **Time-varying multi-seasonal AR models**

Presenter: **Ganna Fagerberg**, Stockholm University, Sweden

Co-authors: Mattias Villani, Robert Kohn

A seasonal AR model with time-varying parameter processes in both the regular and seasonal parameters is proposed. The model is parameterized to guarantee stability at every time point and can accommodate multiple seasonal periods. The time evolution is modeled by dynamic shrinkage processes to allow for long periods of essentially constant parameters, periods of rapid change, and abrupt jumps. A Gibbs sampler is developed with a particle Gibbs update step for the AR parameter trajectories. We show that the near-degeneracy of the model, caused by the dynamic shrinkage processes, makes it challenging to estimate the model by particle methods. To address this, a more robust, faster, and accurate approximate sampler based on the extended Kalman filter is proposed. The model and the numerical effectiveness of the Gibbs sampler are investigated on simulated data. An application to more than a century of monthly US industrial production data shows interesting, clear changes in seasonality over time, particularly during the Great Depression and the recent Covid-19 pandemic.

C1514: **Bayesian nonparametric sensitivity analysis of multiple test procedures under dependence**

Presenter: **George Karabatsos**, University of Illinois-Chicago, United States

A sensitivity analysis method for Multiple Testing Procedures (MTPs) based on marginal p-values is introduced. The method is based on the Dirichlet process (DP) prior distribution, which is specified to support the entire space of MTPs, where each MTP controls either the family-wise error rate (FWER) or the false discovery rate (FDR) under arbitrary dependence among p-values. This DP MTP sensitivity analysis method provides uncertainty quantification for MTPs, by accounting for uncertainty in the selection of such MTPs and their respective threshold decisions regarding which number of smallest p-values are significant discoveries, from a given set of null hypothesis tested, while measuring each p-value's probability of significance over the DP prior predictive distribution of this space of all MTPs, and reducing the possible conservativeness of using only one such MTP for multiple testing. The DP MTP sensitivity analysis method is illustrated through the analysis of over 28,000 p-values arising from hypothesis tests performed on a 2022 dataset of a representative sample of 3 million U.S. high school students, observed across 239 variables. They include tests which, respectively, relate variables about the disruption caused by school closures during the COVID-19 pandemic, with various mathematical cognition, academic achievement, and student background variables. A new R package, bnpMTP, can be used to implement this sensitivity analysis method.

C1515: **Regression on a Lie group: Model and statistical properties**

Presenter: **Johan Aubray**, ENAC, France

Co-authors: Florence Nicol, Anne Francoise Yao

The focus is on the problem of geodesic regression in a Lie group, with the aim of estimating the position of a mobile from noise-measured positions. In the context of air traffic management (ATM), estimating the position of aircraft or a drone can be a complex task, as the data cannot be processed with the usual statistical tools, as in the case of data belonging to a Euclidean space. This regression problem has been previously tackled by modeling the evolution of a mobile (position and orientation) in the Special Euclidean group SE(3). The geodesic that best fits the measurements of a mobile's trajectory in the Riemannian least-squares sense was calculated using the Levi-Civita connection. We study the statistical properties of intrinsic estimators of geodesic regression.

CC400 Room BCB 313 STATISTICAL MODELLING

Chair: Eleonora Arnone

C1088: **Spatio-temporal rainfall forecasting by novel Kumaraswamy-Teissier distribution and seasonal VARMA models**

Presenter: **Kamana Mishra**, Indian institute of technology Mandi, India, India

Co-authors: Neeraj Poonia, Tanmay Kayal, Sarita Azad

A novel three-parameter Kumaraswamy-Teissier distribution is proposed, extending the classical Teissier family, offering enhanced flexibility in modeling skewed data. The key statistical properties of the distribution are derived, including its moments, quantile function, moment generating function, and order statistics. Parameter estimation is performed via maximum likelihood, maximum product spacing, and a Bayesian framework with the Metropolis-Hastings algorithm, which allows posterior inference with credible intervals and highest posterior density regions. Monte Carlo simulations are conducted to compare the efficiency of the estimator. The proposed model is applied to rainfall data from five stations in the Northwest Himalaya and demonstrates a superior fit relative to conventional models. To perform a spatiotemporal analysis, seasonal ARIMA and VARMA models are fitted to KTD-transformed data and raw data. Comparative forecasting reveals that incorporating spatial structure leads to improved predictive performance. Moreover, it is observed that applying these time series models directly to raw rainfall data often yields implausible negative forecasts and prediction intervals. In contrast, applying the VARMA model to KTD-transformed data ensures physical realism and statistical consistency, emphasizing the necessity of proper marginal transformations for modeling. Results highlight the model's potential for applications in hydrology and climate-related studies.

C1433: **Some characterizations of a bivariate Poisson distribution**

Presenter: **Violetta Piperigou**, University of Patras, Greece

Starting from Campbell's bivariate Poisson distribution, various properties can be derived. Certain collections of these properties are demonstrated, than are exclusively possessed by the abovementioned discrete distribution. Moreover, similarities with some established characterization results for the bivariate normal distribution are highlighted.

C1156: **Quantifying model complexity in normalizing flows**

Presenter: **Alexander Ritz**, Clausthal University of Technology, Germany

Co-authors: Benjamin Saefken

Normalizing flows offer a flexible and interpretable way to obtain conditional density estimates based on a transformation model. While there is an intuitive relation between the parameterisation chosen for the underlying transformation function and the complexity of the model, an exact quantification has not been derived so far. Within the literature, there appears to be no clear justification for a particular choice of model other than

the expressiveness of the transformation (universality being preferred usually) and its numerical or analytic convenience. While these are important concerns, deriving a more formalised motivation could help avoid over- and underfitting, allowing researchers to base their model selection on a more objective measure. One way to construct such a model choice criterion may be formulated analogously to established covariance penalties for prediction error estimates, utilizing the underlying normal distribution of the model class. The performance of such an approach is then assessed by means of a simulation study; evaluating the impact of the derived model choice criterion on the prediction error of the chosen model. Likewise, the chosen models' transformation complexity in terms of parameter count and expressiveness is assessed in order to judge their parsimony.

C0443: **Boosting for conditional logistic regression**

Presenter: **Gunther Schauburger**, Technical University of Munich, Germany

Co-authors: Stefanie Klug, Andreas Mayr

Conditional logistic regression or conditional logit models are a standard analysis tool in two distinct research areas, which are the analysis of matched case-control studies and the analysis of discrete choice data. While the models and their parameterization may differ slightly between the research areas, the underlying estimation process is the same. In both cases, using the conditional logit model implies strict assumptions and restrictions for the model. In particular, the covariate effects are related to the response variable in a linear and additive relationship. A general machine learning approach via boosting techniques is proposed, which overcomes these restrictions. It provides the possibility of combining linear, non-linear, tree-based, and spatial effects into one model and allows for data-driven feature selection. A variety of parameterizations can be chosen, which makes it suitable both for matched case-control studies and discrete choice data. For illustration, the method is applied to real-world data. The first application is a matched case-control study on the effect of plasma vitamin E on the first stroke. The second application is about discrete choice data on the choice of travel modes of university students.

CV394 Room Virtual R02 FORECASTING (VIRTUAL)

Chair: Matteo Fontana

C1169: **Pairwise Markov chains for volatility forecasting**

Presenter: **Elie Azeraf**, Azeraf Financial Consulting, France

The pairwise Markov Chain (PMC) is a probabilistic graphical model that extends the classical hidden Markov model (HMM). Despite its flexibility, the PMC has rarely been employed for continuous value prediction, mainly due to challenges in modeling observations within generative frameworks. The aim is to propose a new prediction algorithm for the PMC that overcomes these limitations. The approach (i) resolves the feature modeling problem, fully exploiting the PMCs expressive power, and (ii) provides a general mechanism to extend any predictive model with hidden states that evolve over time, thereby introducing non-stationarity in a principled manner. This methodology is applied to financial volatility forecasting, comparing it with standard benchmarks such as GARCH(1,1) and feedforward neural networks across multiple asset pairs. The empirical results highlight that, under regime changes commonly observed in volatility, the PMC-based extension consistently improves predictive accuracy, demonstrating its practical and theoretical value.

C1253: **Are transformer models good for stock price forecasting?**

Presenter: **Aditya Maheshwari**, Indian Institute of Management Indore, India

Co-authors: Mahima Kumavat, Akshat Vats

Accurate stock price prediction is one of the most challenging tasks in financial decision-making, due to the inherent volatility, nonlinearity, and noise inherent in financial time series. Traditional statistical models are interpretable but fail to capture complex and nonlinear temporal dependencies and long-range patterns. The aim is to present a comparative analysis of state-of-the-art forecasting models, including ARIMA, Long Short-Term Memory (LSTM), Autoformer, Temporal Fusion Transformer (TFT), Informer, PatchTST, and the recently proposed Neural Hierarchical Interpolation for Time Series (NHITS), for five-day-ahead stock price prediction. Using 14 years of historical data of 20 stocks listed on NSE, each model is trained under identical conditions. These models are also validated and evaluated. This ensures that the comparison is fair. The experimental analysis shows that transformer-based models are better at financial time series forecasting. These models are ahead of classical linear approaches in capturing nonlinear dependencies. The temporal fusion transformer (TFT) model achieves the lowest average root mean square error (RMSE) of just 1.41%, showing its superior predictive accuracy. Additionally, PatchTST and NHITS exhibit strong robustness and accuracy, highlighting their effectiveness in modeling local and global patterns using patching and hierarchical interpolation strategies.

C0602: **Words matter: Forecasting economic downside risks with corporate textual data**

Presenter: **Cansu Isler**, Brandeis University, United States

Forecasting downside risks to economic growth is increasingly critical for policymakers and financial institutions. A novel daily sentiment indicator is constructed using textual analysis of corporate disclosures from SEC 10-K and 10-Q filings. Firm-level sentiment is measured as the year-over-year change in the tone of forward-looking statements, based on counts of positive and negative words from the Loughran-McDonald dictionary. These firm-level scores are aggregated into a weekly sentiment index, weighted by market capitalization to reflect broader economic signals. The index is integrated into a mixed data sampling (MIDAS) quantile regression framework to forecast lower quantiles of US GDP growth. Results show that this sentiment-based indicator significantly improves the prediction of economic downturns, outperforming traditional metrics such as the national financial conditions index (NFCI). The approach offers a timely and interpretable measure of market sentiment drawn directly from firms' own assessments of risk and outlook. These findings underscore the potential of corporate textual data as a forward-looking tool for macroeconomic surveillance and policy design.

C1343: **Predicting new sports records: Bootstrap enhancements in extreme value analysis**

Presenter: **Valentina Mameli**, University of Udine, Italy

Co-authors: Michele Lambardi di San Miniato, Federica Giummole, Giovanni Fonseca

A new world record in athletics represents an extreme event, occurring in the tails of the distribution of peak performances. Therefore, extreme value theory provides a natural framework for modeling the distribution of sports records. The commonly used plug-in approach for predicting future records, which replaces unknown parameters with sample estimates in the model, can be unreliable, particularly when sample sizes are small. To address this limitation, two bootstrap-based corrections proposed in the literature are employed: One improves predicted quantiles, while the other, more recent, enhances predicted probabilities, providing more reliable predictions in small-sample contexts. The proposal builds on these two corrections by incorporating the possibility of serial correlation between consecutive annual records in an autoregressive fashion. Applications to records in athletics reveal substantial differences between predictions obtained using the classical plug-in approach and those obtained by the bootstrap methods. These discrepancies highlight the impact of small samples. As a result, the classical method suggests that records are very hard to break, while the bootstrap-based corrections suggest a near future richer in new records.

Sunday 14.12.2025

17:40 - 18:55

Parallel Session J – CFE-CMStatistics 2025

CO064 Room BCB G07 BAYESIAN MIXTURE AND LATENT VARIABLE MODELS FOR LARGE-SCALE PROBLEMS Chair: Mayetri Gupta**C0583: Modelling single-cell RNA-seq data****Presenter:** Indranil Mukhopadhyay, University of Nebraska - Lincoln, United States**Co-authors:** Pronoy Kanti Mondal

Analysis of single-cell RNA-seq data is challenging due to severe sparsity and the influence of many interacting biological factors. Any downstream analysis requires appropriate modeling of raw data, taking care of the inherent complexities present in the data. Bimodal patterns of the expression data add more complexity. The aim is to propose a statistical model that takes care of all these factors simultaneously and fits a probability distribution on the expression levels of each gene, leveraging information from the entire data set. Based on this modeling, a two-sample testing method is also developed for testing differential expression between two groups. Extensive data analysis is performed based on simulated and real data sets to validate the methods.

C0630: A Bayesian hierarchical mixture model for classifying compositional data**Presenter:** Tereza Neocleous, University of Glasgow, United Kingdom**Co-authors:** Catherine Holland, Oliver Stoner

Compositional data refers to multivariate sets of non-negative components, where the primary interest is in the size of the components relative to their total and relative to each other. Such data can be expressed directly as proportions summing to one or measured in absolute terms. Standard statistical techniques are often unsuitable due to the constrained nature of the data, complex correlations, and the presence of zeros, particularly structural zeros, which represent true absence rather than values below detection limits. In forensic glass analysis, compositional data methods have emerged as powerful tools to account for dependencies between elemental components and aid in quantifying the strength of the glass evidence found at crime scenes. The focus is on a forensic elemental glass database that contains a significant number of structural zeros. In these instances, traditional log-ratio transformations, the main technique for compositional data, are undefined. A flexible integrated clustering approach is proposed within a Bayesian hierarchical model for compositional data with structural zeros and a multilevel structure. The model is motivated by the interest in classifying glass fragments by type. Performance was evaluated using five-fold cross-validation, demonstrating superior classification accuracy over less flexible methods. The method offers computational efficiency, supporting the method's practicality for forensic and other real-world applications.

C1273: Located latent class modelling for expert fact-checking from many non-expert checks**Presenter:** Michele Lambardi di San Miniato, University of Udine, Italy**Co-authors:** Michela Battauz, Ruggero Bellio, Paolo Vidoni

Fact-checking is increasingly more critical as false news spreads across social networks. Experts cannot check all the news as fast as needed, so crowd-sourcing is on the rise as a way to distribute tasks to many non-experts. This topic is widely studied in the information retrieval literature. A Bayesian-located latent class model is proposed to surrogate expert judgment by aggregating ratings from multiple workers. The approach combines a prior distribution for the expert's rating with conditional distributions for each worker's rating. Ratings are ordinal variables with levels ranging, e.g., from "false" to "true". Monotonic effects for ordinal predictors and hierarchical priors for workers' parameters help regularize the model. The approach is evaluated on a dataset of fact-checks by PolitiFact, with each statement also rated independently by 10 workers. Patterns in the workers' parameters motivate grouping the workers according to their political orientation. A closed-form but less structured alternative utilizes only the workers' misclassification probabilities and treats all variables as categorical. The hierarchical model outperforms the less expensive alternative and achieves a reasonable misclassification rate.

CO089 Room Virtual R01 RECENT ADVANCES IN FUNCTIONAL DATA ANALYSIS**Chair: Cai Li****C1453: Model-based statistical depth for multivariate sparse functional data****Presenter:** Yue Mu, St. Jude Children's Research Hospital, United States

Functional depth provides a powerful tool for ordering and characterizing complex data, yet most existing methods are developed for univariate functional observations and assume dense sampling. The aim is to introduce a norm-based depth framework tailored for multivariate functional data, where each observation consists of multiple correlated trajectories. By leveraging a reproducing kernel Hilbert space norm, the proposed multivariate depth simultaneously accounts for both within-trajectory variation and cross-component structure, enabling a coherent ranking of multivariate samples. To address the challenges of real-world data, where functional observations are often recorded at irregular and sparse time points, the method is extended to the sparse setting. This allows the norm-based depth to be estimated from limited noisy measurements. Extensive simulation studies are conducted to evaluate the method against existing functional depths. The results demonstrate that the proposed approach not only captures centrality and outlyingness more effectively in multivariate settings, but also preserves reliable performance under sparse designs. Applications to ALS datasets further illustrate the practical utility of the method, showing its ability to detect clinically meaningful outliers and provide consistent characterization of disease trajectories across multiple biomarkers.

C1466: Optimal experimental designs for low-rank function completion**Presenter:** MingHung Kao, Arizona State University, United States

The focus is on optimal experimental designs for collecting high-quality sparse functional data, consisting of observations of one or more random functions $X_1(t), \dots, X_M(t)$ taken at sparse and possibly irregularly spaced points over a compact domain T . Common objectives in analyzing such data include recovering the functions $X_m(t)$ across T , and predicting a response function $Y(t)$ using $X_1(t), \dots, X_M(t)$ as predictors. Drawing on ideas from low-rank matrix completion, a low-rank function completion (LRFC) framework has recently been proposed to efficiently carry out such analyses. However, optimal design strategies to ensure precise inference under the LRFC framework have not been studied, and existing design approaches for sparse functional data analysis (FDA) may become unwieldy under this framework. An efficient method is proposed for identifying optimal designs tailored to the LRFC framework. Its effectiveness is demonstrated, and its applicability to studies where other sparse FDA methods are considered is highlighted.

C1456: Mendelian randomization with pleiotropy through partially functional linear regression**Presenter:** Cai Li, St. Jude Children's Research Hospital, United States

A novel Mendelian randomization (MR) framework is proposed, that models instrumental variables as random functions to account for pleiotropy through a functional partially linear regression. Unlike conventional MR methods that treat SNPs individually, this approach leverages information across entire genes, thereby capturing correlations among SNPs and strengthening signals. The method incorporates a roughness penalty that both respects the structure of functional effects and serves as a regularization to facilitate identifiability of causal effects in the presence of pleiotropy. Building on this, a "smoothness" assumption is introduced, which generalizes the InSIDE (Instrument Strength Independent of Direct Effect) assumption, to further guarantee causal identifiability. For estimation, a penalized partially functional linear regression approach is proposed, implemented as a one-step generalized method of moments (GMM) procedure. This framework enables inference on both causal effects and functional direct effects. The utility of the method is illustrated in uncovering causal relationships between gene expressions and Alzheimer's disease biomarkers. Simulation studies demonstrate that our approach performs favorably compared with state-of-the-art MR methods.

CO023 Room Virtual R02 NEW MACHINE LEARNING AND BAYESIAN TECHNIQUES (VIRTUAL)**Chair: Jairo Fuquene****C0596: A path signature perspective of process data feature extraction****Presenter:** Xueying Tang, University of Arizona, United States

Computer-based interactive items have become prevalent in recent educational assessments. In such items, the entire human-computer interactive process is recorded in a log file as timestamped action sequences. Such response process data are noisy, diverse, and in a nonstandard format. Several methods have been developed to extract information from response processes. However, these methods often focus on the action sequences and ignore the timestamps in response processes. A new feature extraction method is introduced that incorporates the information in both action sequences and timestamp sequences. Based on the concept of path signature, the proposed method extracts features characterizing the response processes at different levels of detail. The proposed method is applied to both simulated data and real response process data from PIAAC to compare the information contained in action and timestamp sequences and demonstrate the potential benefit of incorporating time information for assessing respondents' latent ability.

C1451: Multilevel regression and poststratification**Presenter:** Yajuan Si, University of Michigan, United States

Multilevel regression and poststratification (MRP) is a popular method for addressing selection bias in subgroup estimation, with broad applications across fields from social sciences to public health. The inferential validity of MRP is examined in finite populations, exploring the impact of post-stratification and model specification. To enhance the fitting performance of the outcome model, modeling the inclusion probabilities conditionally on auxiliary variables and incorporating flexible functions of estimated inclusion probabilities as predictors in the mean structure is recommended. A statistical data integration framework is presented that offers robust inferences for probability and nonprobability surveys, addressing various challenges in practical applications. Simulation studies indicate the statistical validity of MRP, which involves a tradeoff between bias and variance, with greater benefits for subgroup estimates with small sample sizes, compared to alternative methods. The methods are applied to the Adolescent Brain Cognitive Development (ABCD) study, which collected information on children across 21 geographic locations in the U.S. to provide national representation, but is subject to selection bias as a nonprobability sample.

C1358: Fast generalized spatial multilevel blockNNGP modelling**Presenter:** Zaida Quiroz, Pontificia Universidad Catolica del Peru, Peru**Co-authors:** Marcos Prates, Dipak Dey, Zhiyong Hu

Typical geostatistical models only consider the case where one response at each location is observed. However, the situation with multiple replicates at spatial locations is seldom discussed. Moreover, the generalized spatial Gaussian process models encounter computational difficulties when the size of the spatial domain becomes massive. Thus, fast generalized spatial multilevel models that use a block nearest neighbor Gaussian process to scale to large datasets are introduced. The proposed method uses integrated nested Laplace approximation (INLA) to avoid long sequential updates of the Markov chain Monte Carlo (MCMC) methods. A simulation study is performed under different response distributions to show the model parameter estimation capacity, computational efficiency, and prediction performance. Finally, the proposed models are fitted to the data of Beijing housing transactions to predict the sales price of houses at unobserved locations. The studies demonstrate that the proposed models have advantages in fitting and prediction, making the interpretation better substantiated.

CO222 Room BCB 206 DYNAMIC PANEL DATA MODELS: A TRIBUTE TO H. PESARAN**Chair: Joachim Schnurbus****C0873: Prediction-oriented modeling of three-dimensional panel data****Presenter:** Markus Fritsch, University of Passau, Germany**Co-authors:** Harry Haupt, Daniel Henderson, Joachim Schnurbus

In the last decade, worldwide migration has been on the rise, and in 2020, an estimated 281 million people (3.6 percent of the world's population) lived in a country in which they were not born. Due to the large proportion of people on the move, modeling international migration flows, understanding their determinants, and predicting expected future migration is of great interest to political decision makers, researchers, and society in general. Starting from a general framework for modeling migration flows based on three-dimensional panel data when accounting for multi-dimensional fixed effects, dynamics, return flows (often referred to as 'reciprocity' in the literature), and asymmetric effects, we detail assumptions frequently encountered in empirical research that lead to a wide range of model specifications. It summarizes how to obtain consistent estimates of the model parameters by providing moment conditions for estimation with the generalized method of moments (GMM) and commenting on inference and the related weighting matrix. The approach is illustrated by modeling decadal and quinquennial migration flows into Europe from 1960 to 2019 based on potentially relevant economic, climate, and conflict variables while accounting for different types of fixed effects.

C0919: Estimation of linear dynamic panel data models based on nonlinear moment conditions**Presenter:** Joachim Schnurbus, University of Passau, Germany**Co-authors:** Markus Fritsch, Andrew Adrian Yu Pua

The focus is on estimating the lag parameter of linear AR(1) panel data models based on nonlinear moment conditions with an instrumental variables (IV) approach. The properties of the estimator under large-n fixed-T settings and under large-n large-T sequential asymptotics are derived. Additionally, the behavior in the near unit root case is investigated, the implications of using an alternative formulation of the moment conditions on the asymptotic results are considered, and the consequences of the presence of predetermined regressors are detailed. Simulation results illustrate the finite sample properties of the estimator compared to alternative estimators based on data-generating processes from the literature. An empirical application demonstrates how to apply the estimator.

C0232: Inference in panel SVARs with cross-sectional dependence of unknown form**Presenter:** Saskia Oetzuerk, Georg-August-Universität Göttingen, Germany**Co-authors:** Lennart Empting, Simone Maxand, Konstantin Wagner

Moving-block bootstrap procedures have become a preferred method to determine the sampling uncertainty of vector autoregressive (VAR) model estimation, most prominently visualized as confidence bands around impulse response functions (IRF). These inferential methods are extended for multivariate time series by the cross-section dimension and compile recursive-design moving-block bootstrap procedures for proxy-identified panel VAR models and their structural IRF. The procedures resample blocks of estimated error terms either in (i) the temporal, (ii) the cross-sectional, or (iii) both dimensions jointly. Their asymptotic assessment and a finite-sample Monte-Carlo study both suggest the preferred use of panel-block resampling in both dimensions when confronted with data properties typically found in empirical panel VAR applications.

CO147 Room BCB 207 NONPARAMETRIC METHODS FOR DEPENDENCE AND STRUCTURAL CHANGE**Chair: B Cooper Boniece****C1143: Inference for quantile change points in high-dimensional time series****Presenter:** Jiaqi Li, University of Chicago, United States**Co-authors:** Likai Chen, Mengyu Xu

Change-point detection methods for quantiles are able to effectively detect structural breaks in extreme values. A novel change-point detection scheme is proposed that utilizes fixed quantiles of moving sums (MOSUM) from high-dimensional time series. The approach employs a MOSUM test statistic that aggregates the component series by the max norm. The asymptotic properties of the proposed test statistic are investigated in

the context of high-dimensional time series, allowing for strong or weak cross-sectional dependence by establishing a powerful uniform Bahadur representation. Specifically, the existing uniform Bahadur representation is extended to the high-dimensional setting for dependent data. To the best of knowledge, this is the first proposal for a change-point detection approach that leverages quantiles from high-dimensional time series. Simulation studies demonstrate the effectiveness of the approach. An application is also presented on a real-world dataset for the value-at-risk in S&P500, which showcases the validity of the method in practical settings.

C1195: Nonparametric changepoint detection with theoretically justified thresholds

Presenter: **Hyeyoung Maeng**, Durham University, United Kingdom

The purpose is to introduce a nonparametric changepoint detection method for univariate data sequences. In the existing literature, the empirical cumulative distribution function (CDF) is often used to detect changes within a nonparametric framework. However, since the empirical CDF depends on a chosen quantile q , the detection power can vary significantly with this choice. To address this issue, some aggregation techniques for building a test statistic have been proposed, although they often lack theoretically justified thresholds. This could possibly lead to failure in controlling the false positive rate, as the limiting distribution of the test statistic under the null hypothesis is not theoretically justified. Two directions are explored to address this issue: 1) filling this gap by proposing a better threshold based on theoretical justification, and 2) proposing a new cost function built on a set of quantiles, whose exact limiting distribution under the null is known.

C1347: Online monitoring for distributional changes with energy distances

Presenter: **B Cooper Boniece**, Drexel University, United States

Co-authors: Lorenzo Trapani, Lajos Horvath

The aim is to propose a nonparametric sequential monitoring procedure for detecting distributional breaks in otherwise stationary time series. The method is based on a class of degenerate U-statistics that includes energy distances and kernel-based maximum mean discrepancy (MMD). The framework accommodates weak dependence and general data types, including multivariate and functional observations. Asymptotic behavior is established under the null, and the distribution of detection delays under alternatives is characterized. Simulations and real-data examples demonstrate strong performance in identifying subtle and mixed-type changes.

CO123 Room BCB 208 QUANTIFYING MODEL SELECTION UNCERTAINTY

Chair: Davide Ferrari

C0213: Measures of uncertainty for model selection

Presenter: **Yuanyuan Li**, Munich Re, United States

Co-authors: Jiming Jiang, Xiaohui Liu

Model selection is a key step in statistical analysis, but its outcome can be sensitive to the data and may not always reflect the true data-generating process. Standard model selection methods often provide a single "best" model, without accounting for the uncertainty in that choice. This can lead to overconfident conclusions and unreliable inference. The aim is to propose two measures to quantify uncertainty in model selection. The first is analogous to a confidence set, identifying a collection of plausible models rather than a single choice. The second focuses on the probability of model selection error. Both methods are applicable to classical settings with a fixed number of candidate models and modern high-dimensional problems. They are conceptually simple, computationally efficient, and supported by both theoretical analysis and empirical results. Their applications are demonstrated by applying them to real-life problems.

C0386: Model selection confidence sets for mixture order selection

Presenter: **Alessandro Casa**, Free University of Bozen-Bolzano, Italy

Co-authors: Davide Ferrari

Determining the number of components in finite Gaussian mixture models is a critical task in clustering and density estimation. Traditional methods based on information criteria often select a single model, potentially overlooking the inherent uncertainty in model selection and resulting in overconfident or inaccurate inferences. To address this, a set-valued estimator, the model selection confidence set, is introduced. This method identifies all mixture orders that are statistically indistinguishable from the best-selected model, using a penalized likelihood ratio screening procedure. The confidence set provides formal coverage guarantees, with a high probability of containing the true number of components. Its width serves as an indicator of data informativeness: a narrower set suggests stronger evidence for a specific order, while a wider set signals greater uncertainty. The method adapts well to the complexity of the data distribution and demonstrates strong performance in simulations across various scenarios. Real data applications further validate its practical advantages and robustness over traditional single-model selection approaches.

C0623: Model selection confidence sets for variable selection in regression via Vuong-type tests

Presenter: **Davide Ferrari**, Free University of Bozen/Bolzano, Italy

A fundamental challenge in regression modeling is selecting the set of relevant predictors from a potentially large collection of candidate variables. Traditional approaches choose a single best model using information criteria or penalized likelihood methods. However, models with different variable subsets often yield similar fits, leading to substantial model selection uncertainty and making it difficult to identify the optimal set of predictors. The model selection confidence set (MSCS) is introduced for variable selection in regression, a set-valued estimator that, with a predefined confidence level, includes the true model across repeated samples. Unlike previous approaches, the method builds these sets using a sequence of Vuong's type tests that accommodate comparisons among both nested and non-nested models. To screen models, each candidate is compared through these tests to a reference model selected by an arbitrary model selection method. Rather than selecting a single model, the MSCS identifies all plausible models that are at least as plausible as this selected reference. Theoretical guarantees are provided for asymptotic coverage, and its practical advantages are demonstrated through simulations and real data analysis.

CO131 Room BCB 209 IMPUTATION TECHNIQUES AND SPATIO-TEMPORAL MODELING

Chair: Abdul-Nasah Soale

C0255: A Jensen-Shannon divergence based k-NN algorithm for missing value imputation in compositional data

Presenter: **Michail Tsagris**, University of Crete, Greece

A novel non-parametric method to impute missing values or rounded zeros in compositional data is suggested. The method is based on the k-NN algorithm, utilizes the Jensen-Shannon divergence, and employs the Frechet mean to allow for more flexibility in the estimation process. As an extra feature, the hyperparameter k can be self-adaptive depending on the pattern of missing values or rounded zeros. Unlike restrictive parametric models, the proposed method makes no assumption about the structure of the data and, most importantly, it is applicable even when compositional data contain structural zero values. Through simulation studies using artificial and real data, the proposed algorithm is superior compared to two competing algorithms at various settings, not only in terms of accuracy but also in terms of computational efficiency.

C0290: Correcting latent class confounder bias in observational studies

Presenter: **Abdul-Nasah Soale**, Case Western Reserve University, United States

Co-authors: Emmanuel Tsyawo

Causal inference is one of the goals of most inferential statistical analysis, especially in business, economics, and health. However, this objective is often unattainable due to potential unmeasured confounding. Bias is addressed in estimating average treatment effects in settings involving latent class confounders induced by proxy and ill-defined categorical covariates, which are correlated with both the treatment and response. A two-step process for debiasing is proposed. In the first step, the latent classes are recovered using model-free sufficient dimension reduction and

clustering techniques. The estimated classes are then incorporated into a mixed model to account for the group structure in the data. The proposed method requires minimal assumptions and yields efficient estimates. An extensive simulation study is included to demonstrate the performance of the proposed method on synthetic data compared to existing methods. Two real applications on medical insurance and energy efficiency are also provided to illustrate the utility of the proposed method in practice.

C0410: Spatiotemporal random field framework for signal detection: Simulation studies and stop-signal task fMRI application

Presenter: **Theophilus Acquah**, University of northern Colorado, United States

The aim is to develop and apply a spatiotemporal random field framework for signal detection in task-based functional magnetic resonance imaging (fMRI). The approach integrates random field theory with functional time series modeling through a voxel-wise two-way repeated measures ANOVA, targeting the dynamic detection of task-related neural activation in a stop-signal task. Using preprocessed fMRI data from 20 participants, activation was assessed within a middle axial brain slice with within-subject factors of task (manual, vocal, pseudoword) and run (Run 1, Run 2). The global test statistic X_{\max} , defined as the maximum Z-statistic across all voxels and time points, enables robust detection of peak activation while controlling for multiple comparisons. F-statistics were transformed into Z-statistics and thresholded conservatively to identify significant regions. Spatial heatmaps and time series plots illustrate how activation profiles vary by task, run, and their interaction. Distinct activation emerged in motor and prefrontal regions, with task, run interactions showing temporally specific effects in frontal areas, supporting models of inhibitory control and adaptive cognitive processing. This framework demonstrates the value of combining random field methods with spatiotemporal inference to uncover complex patterns in fMRI data.

CO047 Room BCB 210 SPATIO-TEMPORAL DEPENDENCE AND COPULA MODELS

Chair: Marta Nai Ruscone

C0641: A copula-based network model for evaluating net bubble risk across asset classes

Presenter: **Giovanni De Luca**, University of Naples Parthenope, Italy

Co-authors: Andrea Montanino

Financial bubbles that periodically form and collapse, manifesting as cycles of sharp growth followed by sudden crashes, pose critical risks to global market stability. As such, timely identification and monitoring of these dynamics are essential for policymakers and financial analysts. The contribution to the literature is the focus on the concurrent emergence of speculative bubbles across different asset sectors, particularly cryptocurrencies and Big Tech equities. Using the backward supremum augmented Dickey-Fuller (BSADF) test, explosive price behaviors and flash crashes are detected and dated. To examine cross-asset dependencies and nonlinear co-movements in net bubble values, a copula-based modeling framework is applied. Additionally, network analysis is conducted to map interconnections during episodes of financial exuberance. Findings reveal that cryptocurrencies exhibit a higher frequency and intensity of speculative episodes compared to Big Tech stocks. However, within the Big Tech universe, Tesla demonstrates bubble characteristics on par with the most volatile cryptocurrencies. The network analysis further underscores a strong interdependence between extreme events in the cryptocurrency sector and Tesla stock behavior, suggesting potential systemic linkages.

C0912: Semi-supervised time series clustering with copulas

Presenter: **Fabrizio Durante**, University of Salento, Italy

Co-authors: Alessia Benevento, Roberta Pappada

Clustering algorithms for time series data play a crucial role as a pre-processing step in building multivariate stochastic models. The focus is on copula-based clustering methods that effectively capture comovements among different time series, independently of their marginal distributions. It begins by reviewing several approaches to clustering random variables using various dissimilarity indices derived from association measures such as Kendall's tau, Spearman's rho, and tail dependence coefficients. Novel algorithms are then introduced that enhance the clustering process by integrating additional deterministic information about the variables within a semi-supervised learning framework. These methods are particularly well-suited for geo-referenced time series, where incorporating spatial context into the dissimilarity measure is essential to uncover meaningful patterns in the data.

C0761: Nested Archimedean copula spatiotemporal approach for hail claim hazard

Presenter: **Melina Mailhot**, Concordia, Canada

Co-authors: Nahid Sadr, Klaus Herrmann

The objective of this project is to model hail claim frequencies for a specific province in Canada. A nested Archimedean copula model is used, with zero-inflated negative binomial GLM margins. Based on theoretical results related to the distortion of copulas, we are able to adapt the spatiotemporal model, which is compared to well-known models, such as INLA and models using spatiotemporal stochastic processes. It turns out the methodology, using distorted copulas, outperforms existing techniques.

CO062 Room BCB 211 ADVANCES IN MULTIVARIATE TIME SERIES

Chair: Alain Hecq

C1027: Broken adaptive ridge in time series regression: The TS-BAR

Presenter: **Orest Prifti**, University of Rome Tor Vergata, Italy

Co-authors: Gianluca Cubadda, Luca Margaritella

The broken adaptive ridge (BAR) estimator is extended to high-dimensional, multivariate autoregressive time series models where the number of variables can exceed the sample size, addressing the curse of dimensionality. The BAR is a sparsity-inducing technique that iteratively reweights the L2-penalty to mimic the selection properties of the non-convex L0 penalization. Additionally, BAR enjoys a grouping effect -where highly correlated covariates are automatically jointly selected- as well as the oracle property, thus performing asymptotically as if the true model were given. In sparse, high-dimensional settings, Monte Carlo simulations show how the BAR can perform better than the LASSO and its variants. Its superior feature selection and forecasting precision are further confirmed by two empirical applications in macroeconomics and finance.

C0964: On mixed causal-noncausal structural VAR models in macro-finance

Presenter: **Lison Christiaens**, Maastricht University, Liege University, Belgium

Co-authors: Alain Hecq, Julien Hambuckers

The aim is to propose a methodology to identify and interpret structural shocks, both causal and noncausal, in mixed causal-noncausal (structural) vector autoregressive (VAR) models. These models are relevant in macrofinance, where noncausal components capture forward-looking behavior, expectation-driven dynamics, and bubble-like episodes in asset prices or policy responses. In the first step, model parameters are estimated using the generalized covariance (GCov) estimator, which relies on nonlinear autocovariances in a non-Gaussian framework. However, it is observed that in multivariate settings, GCov may suffer from estimation instabilities, which is addressed by introducing refinements that enhance robustness. In the second step, identified restrictions are imposed to recover structural shocks and construct impulse response functions (IRFs) that reflect the asymmetric, time-irreversible propagation of causal and noncausal shocks. The methodology is evaluated through Monte Carlo simulations and applied to macroeconomic and financial time series.

C1301: Predictability of funding rates

Presenter: **Emre Inan**, York University, Canada

The purpose is to investigate the out-of-sample predictability of perpetual futures funding rates with a particular focus on Bitcoin contracts traded on Binance and Bybit. Throughout the analysis, one-step-ahead point forecasts are generated from a set of double autoregressive models and evaluated

against standard benchmarks. According to the results, model-based predictions outperform the no-change model both in terms of forecast error and directional accuracy, providing strong evidence for the predictability of the next period's funding rate levels. However, the local analysis shows that the stability of the funding rates is evolving over the evaluation period, suggesting a time-varying degree of their predictability.

CO111 Room BCB 213 ADVANCES IN TIME SERIES ANALYSIS AND FORECASTING
Chair: Tommaso Proietti
C1004: Adaptive forecasting: Implementation (R package forecastADAPT)

Presenter: **Liudas Giraitis**, Queen Mary University of London, United Kingdom

Co-authors: Violetta Dalla, George Kapetanios

Forecasting strategies that are robust to structural breaks and structural changes have earned renewed attention in the literature. These strategies are particularly attractive to applied econometricians because they are built on weighted averages of down-weighted past information. They are easy to implement, and in general, they are robust to different types of structural change: they are capable of adapting to it in real time. The aim is to introduce the R package forecastADAPT that implements the robust forecasting strategies via a set of examples and applications that might be of interest to applied econometricians. The theory behind this R package, how the forecast is constructed, the forecast error estimated, and the data-based tuning parameter selected, is based on a prior study.

C1046: Frequency singular spectrum analysis

Presenter: **Pilar Poncela**, Universidad Autonoma de Madrid, Spain

Co-authors: Gabriel Martos, Diego Fresoli

Singular spectrum analysis (SSA) is a nonparametric method for time series modeling and forecasting. Via the singular value decomposition on the so-called trajectory matrix, or equivalently, by diagonalizing the second moment matrix of the data, SSA decomposes a time series into quasi-orthogonal components that aim to maximize variance. The resulting trendlines provide natural estimates of the underlying unobserved trend, cycles, and noise. However, these trendlines are not directly associated with specific oscillation frequencies. A novel extension of SSA that reconciles frequency identification with variance-based decomposition is introduced. The aim is to propose a consistent estimator of the spectral density embedded within the method and offer guidance on selecting appropriate decomposition parameters. Additionally, inferential tools are developed to address the grouping problem, enabling the identification of statistically significant components related to specific oscillation frequencies. The performance of the proposed methodology is demonstrated through simulation studies. The method is further applied to analyze US temperature dynamics, identifying a steady increase of 5F in average temperature over the past century.

C1373: Robust CDF filtering of a location parameter

Presenter: **Alessandra Luati**, Imperial College London, United Kingdom

Co-authors: Leopoldo Catania, Andrew Harvey

The purpose is to introduce a novel framework for designing robust filters associated with signal plus noise models having symmetric observation density. The filters are obtained by a recursion where the innovation term is a transform of the cumulative distribution function of the residuals. The latter downweights extreme values by construction and allows the filters to be analytically tractable. The updating scheme naturally arises as the solution of an optimization problem, where the objective function is a continuous version of the quantile check function, formerly employed as a proper scoring function for quantiles and used to construct robust minimum contrast estimators. Stationarity, ergodicity, and invertibility are derived under minimal assumptions and preserved under different parametric specifications. Estimation is carried out by the method of maximum likelihood, and the asymptotic theory is developed under misspecification. As an illustration, the new filters are applied to brain scan data and compared across Gaussian, Student t, Cauchy, and Logistic density specifications, with alternative methods. Additional results include a novel class of score-driven models and a sub-Gaussian density suitable for robust filtering and modelling, arising as the infinite sum of independent non identically distributed uniform random variables.

CO233 Room BCB M202 TOPICS IN FINANCE AND ECONOMETRICS
Chair: Florian Richard
C1142: Equity market-neutral strategies using variable selection and regularized regression

Presenter: **Federico Severino**, Universite Laval, Canada

Co-authors: Marzia Cremona, Charles-Edouard Sarault

Equity market-neutral strategies are designed to feature no exposure to market risk. The implementation of such strategies relies upon the estimation of the beta coefficients. Unfortunately, traditional beta estimation methods used to implement these strategies often suffer from weak out-of-sample performance, leading to suboptimal ex-post neutrality. Therefore, it is explored how machine learning techniques, particularly variable selection and regularized regression methods, can address the issues faced in the traditional development of equity market-neutral strategies. A range of methods is tested, including ridge regression, lasso regression, and stepwise regressions, to construct portfolios that achieve better ex-post neutrality and lower transaction costs when compared to strategies based on traditional multiple regression models. The results demonstrate that the tested techniques enhance portfolio performance and help minimize trading costs by selecting an optimal number of risk factors to hedge. Research contributes to the academic literature on machine learning in asset pricing and offers practical insights for portfolio management.

C1435: Simultaneous inference in possibly incomplete multivariate regressions with application to asset pricing

Presenter: **Lynda Khalaf**, Carleton, Canada

Co-authors: Marie-Claude Beaulieu, Jean-Marie Dufour

The focus is on finite sample tests and confidence procedures for a system of several estimating equations with cross-dependent, possibly non-Gaussian disturbances, and allowing for possibly missing explanatory variables. Standard statistics, based on the roots of a determinantal equation involving the restricted and unrestricted least squares residuals, are considered. First, the exact distribution of these roots are characterized, and pivotal special cases are discussed. Second, an identifiable ex-post parameter is introduced that embeds the impact of left-out explanatory variables, in which case usual tests remain valid from a conservative perspective. Third, exact simultaneous confidence bands are derived for a vector of regression coefficients, or for a combination thereof. Fourth contribution pertains to asset pricing, with a focus on risk-return analysis. The distinction between the alpha (risk-adjusted outperformance) and beta (systematic risk component) of returns is revisited, from the ex-post parametrization perspective. Empirically, the risk-return profile of catastrophe bond mutual funds is analyzed.

C1284: Insider trading reporting, market activity, and liquidity

Presenter: **Florian Richard**, Universite Laval, Canada

How does insider trading affect volume and liquidity? Using a panel dataset of over 200,000 insider transactions, it is found that the publication of an insider transaction report is associated with a 3% increase in same-day trading volume, followed by a 20-day decline. Purchases are associated with higher trading volume, while sales are associated with lower trading volume over time. Overall, insider transactions widen same-day bid-ask spreads. Market participants are more active in response to larger sales. Transaction size has no effect on spreads. The use of novel spread measures and panel data offers contrasting results with the previous literature.

CO162 Room BCB 307 RECENT ADVANCES IN HIDDEN MARKOV MODELS**Chair: Beatrice Foroni****C0672: Capturing static and dynamic heterogeneity in a Hidden Markov model for binary data****Presenter:** Dalila Failli, University of Perugia, Italy**Co-authors:** Maria Francesca Marino, Francesca Martella

A model-based clustering approach for multivariate binary longitudinal data that accounts for both time-varying and time-constant sources of unobserved heterogeneity. Specifically, a latent Markov model (LMM) specification is considered to capture the dynamic nature of the data and enable dynamic clustering of units over time. Additionally, a multidimensional continuous latent trait is incorporated to capture residual time-constant heterogeneity among units within the same latent state at any given time point, in terms of their responses to the multivariate binary variables. For parameter estimation, the standard Baum-Welch algorithm is extended to accommodate the presence of the continuous latent trait. In this context, the solution of multidimensional integrals not available in closed form is required. A variational approximation is considered to overcome the issue. The performance of the proposed approach is evaluated in terms of parameter recovery and clustering accuracy. The ability of different model selection procedures in identifying the optimal number of latent states and latent trait dimensions is also evaluated.

C0721: Sparse estimation in Markov regime-switching models**Presenter:** Abbas Khalili, McGill University, Canada**Co-authors:** Gilberto Chavez Martinez

Markov regime-switching models are widely used to model heterogeneous and complex relationships in multivariate time series. While maximum likelihood estimation (MLE) is the standard approach for parameter estimation in these models, it often becomes unstable even in moderate parameter space dimensions. A class of regularization-based estimators, designed to address this challenge, is presented. Both the theoretical properties and finite-sample performance of the proposed methods are discussed, followed by an application to real data.

C1107: Non-homogeneous Markov-switching generalized additive models for location, scale, and shape**Presenter:** Timo Adam, Bielefeld University, Germany**Co-authors:** Katharina Ammann

The aim is to propose an extension of Markov-switching generalized additive models for location, scale, and shape (GAMLSS) that allows covariates to influence both the parameters of the state-dependent distributions and the transition probabilities. Traditional Markov-switching GAMLSS combines distributional regression with latent-state time series modelling, but typically assumes constant transition probabilities, which prevents regime shifts from responding to covariate-driven changes. The approach overcomes this limitation by allowing the transition probabilities to vary with covariates, thereby capturing covariate-dependent regime dynamics. The proposed methodology is evaluated through simulation studies, and its practical usefulness is illustrated in a case study on Spanish energy prices.

CO173 Room BCB 308 STATISTICAL INNOVATIONS FOR DEPENDENT DATA**Chair: Hossein Moradi Rekabdararkolaee****C0875: Bayesian closure modeling for dynamic systems****Presenter:** Toryn Schafer, Texas A&M University, United States

Closure modeling is a key challenge in the simulation of dynamical systems, where unresolved processes must be represented accurately to ensure predictive fidelity. A Bayesian approach to closure modeling is explored by casting the problem as variable selection over a latent derivative process governed by an ordinary differential equation. Building on recent work in dynamic discovery, the focus is on scalable computational strategies for posterior inference, including implementations that target GPU architectures. This framework enables rigorous uncertainty quantification while maintaining computational tractability, even in the presence of sparse or noisy observations. The method is demonstrated on an ecological case study, and broader implications are discussed for dynamic model discovery in scientific domains.

C0887: Changepoint detection in categorical time series with application to daily total cloud cover in Canada**Presenter:** Mo Li, UL at Lafayette, United States**Co-authors:** Qiqi Lu, Xiaolan Wang

Changepoints are essential for homogenizing categorical time series and analyzing their trends and variations. The original total cloud cover in Canada was recorded hourly in tenths (or eighths), exhibiting inherent seasonality and serial correlation. A prior study introduced an extended cumulative logit model to detect shifts in the annual frequencies of cloud cover conditions. While annual aggregation mitigates seasonality and serial correlation, it shortens the time series and may lead to overdispersion. A marginalized transition model is introduced to detect a single changepoint in periodic and serially correlated categorical time series. The model captures serial dependence using a first-order Markov chain and enables category-specific changepoint specification. To enhance computational efficiency, a new parameter estimation procedure is developed for obtaining maximum likelihood estimates. A maximally selected likelihood ratio test statistic is then proposed to test for sudden changes in categorical time series, and the method is illustrated using daily total cloud cover observations recorded at 9 a.m. and 3 p.m. at Fort St. John Airport, British Columbia, Canada.

C1140: Modeling spatio-temporal extremes via conditional variational autoencoders**Presenter:** Likun Zhang, University of Missouri, United States

Extreme weather events are widely studied in the fields of agriculture, ecology, and meteorology. Enhanced scientific comprehension of the spatio-temporal dynamics of these events could significantly improve policy formulation and decision-making within these domains. In this paper, we propose a novel approach to model spatio-temporal extremes by integrating climate indices using conditional variational autoencoders (extreme-CVAE). The alignment between modeled and true extremal dependence structures showcases the model's ability to be a spatio-temporal extreme emulator. Along with the decoding path, a convolutional neural network was built to investigate the relationship between climatological dynamics and latent-space parameters, thereby inheriting the underlying temporal dependence structures. The extensive simulation validated the effectiveness and time efficiency of the proposed model. Furthermore, we apply our method to analyze monthly maximum Fire Weather Index (FWI) values over eastern Australia from 2014 to 2024, using GEOS-5 data from the Global Fire Weather Database (GFWED). This case study highlights the model's practical utility and performance in real-world scenarios.

CO316 Room BCB 309 RECENT ADVANCES IN MATRIX AND TENSOR TIME SERIES MODELS**Chair: Jingyang Li****C1290: Multilevel main effects matrix factor model****Presenter:** Clifford Lam, London School of Economics and Political Science, United Kingdom**Co-authors:** Zetai Cen, Kaixin Liu

The aim is to propose a multilevel main effects matrix factor model (MMEFM) to extract meaningful row and column effects from groups of matrix data. As a generalization of a multilevel matrix factor model, rigorous definitions of MMEFM and identifications are provided, together with iterative algorithms for the estimators of all global and local components in the main effects and common component, including global/local main effects and core ranks, with rates of convergence spelt out. An important feature of the model is that all main effects can be nonstationary and still be consistently estimated at each time point. The objective is to propose a test for testing if a multilevel matrix factor model is sufficient against the alternative that MMEFM is necessary for the data. The power of the test, together with the accuracy of various estimators, is demonstrated in

extensive simulation settings. Estimation performance is also compared to other state-of-the-art methods for matrix time series using simulations and a real data set.

C1302: **Dynamic matrix factor model for counts data**

Presenter: **Han Xiao**, Rutgers University, United States

A dynamic factor model is considered for matrix time series, where the observations are counts data. The model is formulated as Poisson observations conditional on the rate matrices, which have log-normal distributions. The logarithm of the rate matrices has the form of a dynamic Gaussian factor model of matrix time series, where the dynamics are captured by the factor process. A log moment method is proposed to estimate the loading matrices, and use the variational inference to estimate the factor process. An autoregressive model is imposed on the factor process to enable predictions. It is also considered to include a trend component in the log rate matrices to capture possible nonstationarity. Theoretical and numerical analyses are conducted for the proposed model.

C1401: **Local interaction autoregressive model for high dimension time series data**

Presenter: **Jingyang Li**, University of Michigan, United States

Co-authors: Yang Chen

High-dimensional matrix- and tensor-valued time series arise in fields such as economics, geophysics, and environmental science. Traditional vector autoregressive models are infeasible in these settings due to excessive parameters and loss of spatial structure under vectorization. Existing matrix autoregressive (MAR) models also face two limitations: They assume each location depends on the entire spatial field from the previous step, while in practice, local interactions dominate; and they impose restrictive low-dimensional structures on the coefficient matrices, limiting flexibility in capturing heterogeneous dependencies. A local interaction autoregressive (LIAR) model is proposed that incorporates spatial locality into matrix and tensor autoregression. Each entry depends only on a neighborhood in past observations, with neighborhoods allowed to vary across locations. A parallel least squares estimator is developed with closed-form solutions for efficient large-scale computation. To further reduce complexity, a separable variant, SP-LIAR, is introduced which preserves flexibility with fewer parameters. Asymptotic theory is established, including consistency, asymptotic normality, and consistent neighborhood selection under broad regimes. Simulations and real data applications show that LIAR offers an effective balance of interpretability, flexibility, and computational efficiency.

C1529: **Modewise additive factor model for matrix time series**

Presenter: **Yuefeng Han**, University of Notre Dame, United States

Matrix-valued time series arise in domains like finance, retail analytics and environmental science, where capturing structured dynamics along both rows and columns is essential. We propose a modewise additive factor model for matrix-valued time series that separates latent dynamics along rows and columns, enabling, for example, store-level and product-level trend analysis. Under mild conditions, the loading spaces for both modes are identifiable up to orthogonal rotation, without requiring restrictive covariance or independence assumptions. To recover these spaces, we develop two estimation algorithms, Modewise INner-product Eigendecomposition (MINE) and COMplement-Projected Alternating Subspace Estimation (COMPAS), and establish their near optimal convergence rates under standard mixing and tail probability conditions. Through simulations and a retail sales case study, we demonstrate the methods accuracy, efficiency and interpretability, offering a versatile tool for structured temporal data analysis.

CO125 Room BCB 310 RECENT APPROACHES TO ENVIRONMENTAL AND SPATIO-TEMPORAL STATISTICS

Chair: Tiffany Tang

C0209: **Neural classification of asymptotic (in)dependence**

Presenter: **Troy Wixson**, University of Massachusetts Amherst, United States

Co-authors: Daniel Cooley

Studies in extremes often aim to extrapolate into the tail beyond the range of the data; For example, to assess the risk of the combined effect of extreme precipitation and storm surge. Extrapolation under the wrong dependence regime can have large negative effects, and thus classification is a necessary early choice in the modeling of multivariate extremes. Inference about the dependence regime is complicated as the regimes are defined asymptotically. A series of experiments is performed to determine whether a finite sample has enough information for a convolutional neural network to reliably distinguish between these regimes in the bivariate case. A new classification tool is developed for practitioners, which is called *nnadic*, as it is a neural network for asymptotic dependence/independence classification. This tool accurately classifies over 97% of test datasets and is robust to a wide range of sample sizes. The datasets that are unable to be correctly classified tend to either be nearly exactly independent or exhibit near-perfect dependence, which are boundary cases for both the ADep and AInd models used for training. These experiments highlight that ADep and AInd models do not so much differ in the strength of tail dependence they can capture (as both regimes can range from independence to complete dependence), but they instead differ in whether the dependence completely decays in the limit, irrespective of the path of that decay.

C0208: **Spatial hyperspheric models for compositional data**

Presenter: **Michael Schwob**, Virginia Tech, United States

Compositional data are an increasingly prevalent data source in spatial statistics. Analysis of such data is typically done on log-ratio transformations or via Dirichlet regression. However, these approaches often make unnecessarily strong assumptions (e.g., strictly positive components, exclusively negative correlations). An alternative approach uses square-root transformed compositions and directional distributions. Such distributions naturally allow for zero-valued components and positive correlations, yet they may include support outside the non-negative orthant and are not generative for compositional data. To overcome this challenge, the elliptically symmetric angular Gaussian (ESAG) distribution is truncated to the non-negative orthant. Additionally, a spatial hyperspheric regression model is proposed that contains fixed and random multivariate spatial effects. The proposed model also contains a term that can be used to propagate uncertainty that may arise from precursory stochastic models (i.e., machine learning classification). The model is used in a simulation study and for a spatial analysis of classified bioacoustic signals of the Dryobates pubescens (downy woodpecker).

C0897: **Gaussian processes and spatial data fusion: Avoiding numerical integration**

Presenter: **Lucas da Cunha Godoy**, University of California Santa Cruz, United States

Co-authors: Marcos Prates, Fernando Quintana, Jun Yan

Spatial data fusion (SDF) combines data on a single phenomenon collected at different spatial resolutions, such as point-referenced measurements from monitoring stations and areal data from satellites or computer models. Standard models assume both data types are realizations of a common Gaussian process (GP), with real data defined as spatial averages of this process. This approach, however, lacks a closed-form solution and requires numerical integration for approximation. This step is not only computationally expensive but also relies on an arbitrary grid selection, for which there is no consensus on an optimal resolution. As an alternative, a generalization of the isotropic GP used in SDF is proposed. The method defines the covariance function based on a distance between sets, allowing for a direct calculation of covariance between any combination of points and areas. This approach completely bypasses the need for integration. The theoretical conditions are established that ensure the validity of the model, and the method's utility is illustrated using an atmospheric temperature dataset. The proposed methodology offers a computationally efficient and conceptually simple alternative, removing the ambiguity of grid selection and providing a more robust solution for practitioners.

CO201 Room BCB 311 BRIDGING CAUSALITY, CONNECTIVITY, AND ROBUST INFERENCE**Chair: Ashkan Ertefaie****C0601: Resilience measures for the surrogate paradox****Presenter:** Layla Parast, University of Texas at Austin, United States

Surrogate markers are often used in clinical trials to evaluate treatment effects when primary outcomes are costly, invasive, or take a long time to observe. However, reliance on surrogates can lead to the "surrogate paradox, where a treatment appears beneficial based on the surrogate but is actually harmful with respect to the primary outcome. Formal measures are proposed to assess resilience against the surrogate paradox. The setting assumes an existing study in which the surrogate marker and primary outcome have been measured (Study A) and a new study (Study B) in which only the surrogate is measured. Rather than assuming transportability of the conditional mean functions across studies, a class of functions is considered for Study B that deviates from those in Study A. Using these, the distribution of potential treatment effects is estimated on the unmeasured primary outcome and defines resilience measures, including a resilience probability, resilience bound, and resilience set. The approach complements traditional surrogate validation methods by quantifying the plausibility of the surrogate paradox under controlled deviations from what is known from Study A. The performance of the proposed measures is investigated via a simulation study and application to two distinct HIV clinical trials.

C0911: What is in a stimulus: Exploring functional connectivity and causal relationships**Presenter:** Reza Ramezan, University of Waterloo, Canada

Neurons transmit information through consecutive electrochemical waves, which are also called spike trains due to their temporal localization. Modelling spike trains falls naturally within the statistical point process framework. However, developing multivariate point process models for neural spike trains that are flexible, computationally efficient, and biologically meaningful remains challenging. The multivariate Skellam process with resetting (MSPR) and its continuous-space generalizations within a latent factor model framework are some options to address these challenges. While these models, like most alternatives, effectively capture functional connectivity between neurons, causal relationships in synaptic transmission have been less explored. A key factor influencing functional connectivity in neuronal ensembles is the stimulus signal to which neurons respond. Using experimental data from the visual cortex of rhesus monkeys, the causal effects of stimulus signals on functional connectivity are analyzed, alongside comparisons of general neuronal associations accounting for stimulus influence. The findings provide insight into the dynamics of neural communications under external stimulation.

C0974: Stochastic interventions, sensitivity analysis, and optimal transport**Presenter:** Alexander Levis, University of Pennsylvania, United States**Co-authors:** Edward Kennedy, Alexander McClean, Sivaraman Balakrishnan, Larry Wasserman

Recent methodological research in causal inference has focused on effects of stochastic interventions, which assign treatment randomly, often according to subject-specific covariates. It is demonstrated that the usual notion of stochastic interventions has a surprising property: When there is unmeasured confounding, bounds on their effects do not collapse when the policy approaches the observational regime. As an alternative, the purpose is to study generalized policies, treatment rules that can depend on covariates, the natural value of treatment, and auxiliary randomness. It is shown that certain generalized policy formulations can resolve the "non-collapsing" bound issue: Bounds narrow to a point when the target treatment distribution approaches that in the observed data. Moreover, drawing connections to the theory of optimal transport, generalized policies are characterized that minimize worst-case bound width in various sensitivity analysis models, as well as corresponding sharp bounds on their causal effects. These optimal policies are new and can have a more parsimonious interpretation compared to their usual stochastic policy analogues. Finally, flexible, efficient, and robust estimators are developed for the sharp nonparametric bounds that emerge from the framework.

C1501: A framework for mark scaling**Presenter:** Benjamin Baer, University of St Andrews, United Kingdom

Sometimes, the raw mark/grade of a student in a module/course may not be deemed suitable, so it will get scaled/curved to an adjusted mark. After reviewing some existing scaling methods, we state some desiderata for the scaling and study whether existing methods meet the desiderata. Finally, we present a refined method.

CO200 Room BCB 312 INFERENCE AND TESTING IN COMPLEX BIOMEDICAL STUDIES**Chair: Florin Vaida****C1316: Generalized estimating equation for cell-cell correlation in single-cell RNA seq data****Presenter:** Xinlian Zhang, University of California, San Diego, United States**Co-authors:** Toni Gui, Tuo Lin, Xin Tu

For analyzing the single-cell RNA sequencing data, it is believed that cells from the same individual share common genetic and environmental backgrounds and are not statistically independent. Many popular methods, such as the default Wilcox test in the FindMarker function in the Seurat package, do not address this issue, leading to potentially biased inference. There are more recent works arguing for the generalized linear mixed models with a random intercept for the individual, to properly account for the correlation among measures from cells within an individual. However, the traditional mixed effect model has strong assumptions requiring the same and strictly positive correlation across all cells in one individual. It is demonstrated that this can be rather restrictive for real data seen, given the strong heterogeneous nature of all cells. In the case of a violated positive correlation assumption, the classical random effects model demonstrates consistent biased inference. The aim is to propose the usage of the generalized estimating equation-based semi-parametric approach for this issue and demonstrate its robust and efficient performance in both simulation and real data that focuses on revealing common and unique gene expression signatures in primary CD4+ T cells latently infected with HIV under different conditions.

C1457: Issues in causal inference using matching with longitudinal survey data**Presenter:** Karen Messer, University of California, San Diego, United States**Co-authors:** Natalie Quach, Chen Jiayu

Matching methods are well-established as a popular approach to inference, with well-developed tools and methods, and effective use in the public health literature. Several problems encountered in population studies of smoking behavior using longitudinal survey data are outlined, including questions of mediation and of adjusting for prognostic variables. In the setting of large-scale complex surveys, additional questions always include how and when to incorporate survey weights, and how to carry out resampling-based approaches to inference.

C1489: Modeling events with competing causes of known exposure: Risk of concussion due to sports in children**Presenter:** Florin Vaida, University of California San Diego, United States**Co-authors:** Wenjing Meng

In the ongoing ABCD study, over 10,000 children at the age of 9-10 are followed up for up to 10 years, with extensive information collected at baseline and yearly. We are concerned with the risk of mild traumatic brain injury (mTBI), which may be caused by playing sports, as reflected in the baseline study visit. About 4% of children report mTBI at baseline, but the cause of concussion is not reported. Information about sports played in the past, and the exposure, in the form of average time per week, is recorded. We model the risk of mTBI as a function of individual sports played, as well as other demographics and socio-economic characteristics of the children and families, using an extension of Poisson regression, via a mixture of censored Poisson distributions, with censoring at 1 event. The exposure determines the mixing probabilities. The maximum likelihood estimator is computed using the EM algorithm, with standard errors determined via the Louis formula. We show that the hazard of mTBI is more

than 10 times when playing a sport such as soccer (football), compared to time spent not playing sports.

CO257 Room BCB 313 RECENT ADVANCES IN STATISTICAL LEARNING ON COMPLEX DATA SETS
Chair: Tianxi Li
C0331: Robust multi-ancestry PWAS utilizing Bayesian fine-mapping

Presenter: Haoran Xue, City University of Hong Kong, Hong Kong

Co-authors: Chengli Zhang, Chong Wu

Proteome-wide association studies (PWAS) have emerged as a powerful tool for identifying proteins associated with complex diseases, which can serve as potential drug targets. The conventional PWAS approach employs a two-stage least squares (2SLS) regression, utilizing genetic variants as instrumental variables (IVs). However, the validity of this approach can be compromised by the widespread pleiotropy of genetic variants, which can lead to the identification of false-positive causal proteins for diseases. Furthermore, the varying linkage disequilibrium (LD) patterns and effect sizes of genetic variants across ancestries limit the power of PWAS when analyzing different populations separately. To address these challenges, a robust and powerful PWAS method is proposed that integrates proteomics and disease data from multiple ancestries and employs a Bayesian fine-mapping approach to detect invalid IVs. The effectiveness of the proposed method is demonstrated by applying it to a large-scale biobank dataset, identifying putative causal proteins for complex human diseases.

C0989: Causal learning with invalid instruments for high-dimensional imaging responses

Presenter: Shan Yu, University of Virginia, United States

Co-authors: Yiting Wang, Chunlin Li

Large-scale neuroimaging studies provide high-dimensional imaging phenotypes for understanding brain-related diseases and their associations with genetic, environmental, and clinical factors. However, causal inference remains difficult due to unmeasured confounding. Instrumental variables (IVs) are commonly used to address confounding, but identifying valid IVs is often difficult in practice. Existing IV methods face significant limitations when applied to high-dimensional imaging data, particularly in handling large-scale data and incorporating the intrinsic spatial structure of images. To address these challenges, a novel framework is proposed, causal image-on-scalar regression with invalid instrumental variables (CISR-IIV), which enables estimation of spatially varying causal effects in the presence of potentially invalid IVs. The approach integrates nonparametric spatial smoothing to capture the spatial structure of imaging data, combined with a lasso-based instrumental variable selection strategy to handle potentially invalid instruments. Rigorous theoretical guarantees are established for the CISR-IIV framework, including selection consistency and the asymptotic distribution of the estimated causal effects. Building on these theoretical results, asymptotic confidence intervals and data-driven simultaneous confidence regions are further constructed. CISR-IIV is applied to the Alzheimer's disease neuroimaging initiative to illustrate its effectiveness.

C1438: Streaming tensor decomposition in imaging analysis

Presenter: Xiwei Tang, University of Texas at Dallas, United States

Co-authors: Haowen Zhou

The analysis of complex medical imaging data, particularly from longitudinal or streaming settings, presents significant challenges due to the dynamic and evolving nature of pathological changes. For example, magnetic resonance imaging (MRI) and diffusion imaging (DI) provide critical insights into brain microstructure, but capturing and analyzing the temporal patterns within such multidimensional data streams is a complex task. Tensor decomposition has shown promise in a variety of applications, from neuroscience to social networks, for extracting meaningful patterns. The aim is to present a novel tensor-based framework that combines advanced streaming tensor decomposition techniques with the analysis of longitudinal imaging data. The method is designed to track temporal changes and identify regions that exhibit longitudinal patterns. It is built on tensor tracking, incorporating streaming tensor decompositions to handle the complex, real-time data streams generated by longitudinal imaging. By addressing the unique challenges of temporal structure analysis and data stream factorization, this approach provides new insights into the progression of medical imaging analysis.

CO135 Room BCB 402 MODERN STATISTICAL FRONTIERS FOR BIOMEDICAL AND HIGH-DIMENSIONAL DATA
Chair: Matteo Borrotti
C1315: Penalized deep partially linear cox models with application to CT scans of lung cancer patients

Presenter: Yuming Sun, College of William and Mary, United States

Lung cancer is a leading cause of cancer-related death worldwide, underscoring the need to understand mortality risks for effective, personalized treatment. The National Lung Screening Trial (NLST) used CT texture analysis to quantify image-based risk factors. Partially linear Cox models are increasingly used in survival analysis for integrating both traditional (e.g., age, clinical covariates) and novel (e.g., imaging features) risk factors by combining parametric and nonparametric components. However, when the number of parametric covariates exceeds the sample size, model fitting becomes challenging, and nonparametric modeling suffers from the curse of dimensionality. A novel penalized deep partially linear Cox model (Penalized DPLC) is proposed, which integrates the SCAD penalty for selecting important texture features and uses a deep neural network to estimate the nonparametric component. Theoretical guarantees are established for the estimator, and its superior performance is demonstrated through simulations in both prediction and feature selection. Finally, the method is applied to NLST data to investigate how clinical and imaging features relate to survival, offering insights into the prognostic value of these multimodal risk factors.

C1319: Advances in microbiome differential abundance analysis: Group-wise normalization and multivariate count regression

Presenter: Kyu Ha Lee, Harvard T.H. Chan School of Public Health, United States

The focus is on a new perspective on normalization for differential abundance analysis, highlighting a group-level approach to mitigate compositional bias. Within this framework, two novel strategies are introduced, group-wise relative log expression and fold-truncated sum scaling, that enhance statistical power while maintaining false discovery control. Then, recent developments in multivariate regression modeling for zero-inflated count data are discussed, with a particular focus on Bayesian variable selection. These models enable simultaneous inference across multiple taxa, scale effectively to high-dimensional datasets, automatically identify key associations, and integrate normalization procedures that address compositional challenges. Simulation studies and real-data applications illustrate that this combined framework surpasses existing univariate and multivariate methods, offering a robust and versatile toolkit for microbiome data analysis.

C1434: Mitigating bias in analyzing privacy-preserved time-to-event data

Presenter: Yi Xiong, University at Buffalo, United States

Sharing time-to-event data is beneficial for enabling collaborative research efforts, facilitating the design of effective interventions, and advancing patient care. However, sharing the exact survival curves poses concerns over privacy. Although there are several popular privacy-protecting solutions (e.g., binning, differential privacy) that offer strong protection on the data, the "sanitized" data usually has low utility and can result in misleading statistical inference. The aim is to investigate the distortion in bias and variance in regression analysis of sanitized survival data under popular privacy-protecting solutions, and to provide a strategy to mitigate the bias in estimators with sanitized survival data.

C1531: Uncovering high-dimensional genomic signals modulating biological networks via Gaussian graphical models

Presenter: Samuel Anyaso-Samuel, National Cancer Institute, United States

Gaussian graphical models (GGMs) are powerful tools for characterizing direct relationships among biological traits (e.g., gene expression, protein, microbial taxa) through partial correlation coefficients (PCCs). While traditional applications of GGMs focus on static biological networks, the

influence of genomic factors (e.g., SNPs, mutations) on these networks remains underexplored. We propose a two-stage penalized regression framework to identify genomic regulators that modify trait-trait (edges) relationships. Our strategy first constructs a baseline network and then tests for genomic modifiers only among selected trait pairs, reducing computational and multiple testing burdens. Moreover, we develop an analytic procedure to control the false discovery rate (FDR) in regularized regressions, making large-scale analysis feasible. Through simulations, we demonstrate its computational efficiency, accurate Type-I error control, and high sensitivity. We demonstrate the utility of our framework in three case studies: identifying host genetic variants associated with oral microbiome networks, linking gut microbiome taxa to metabolite networks, and uncovering somatic mutations influencing gene expression networks in lung adenocarcinoma.

CO150 Room BCB 403 STATISTICAL MODELING FOR BIOMEDICAL AND PUBLIC HEALTH APPLICATIONS Chair: Dhrubajyoti Ghosh

C0332: The SIR model under missing data: A marginal likelihood approach via dynamical survival analysis

Presenter: **Suchismita Roy**, Duke University, United States

Co-authors: Jason Xu, Alexander Fisher

The SIR model is widely used for modeling epidemic dynamics. Despite its widespread use, parameter inference in the presence of missing data is challenging due to the intractability of the likelihood in such compartmental models. To address this, a closed-form likelihood is developed for incidence data using the dynamical survival analysis (DSA) method, which provides a survival analysis-based interpretation of the SIR model. The method is flexible and computationally efficient. Through simulation, its performance is compared in parameter estimation with other methods that rely on the exact posterior of the SIR model. To further demonstrate the adaptability of the approach, the likelihood is extended to frailty models, illustrating how it can be modified to incorporate individual heterogeneity. Finally, the method is applied to real-world data, demonstrating its practical utility for epidemic inference with limited observations.

C0881: InterSpatial: Leveraging low-resolution spatial transcriptomics to infer cell-cell communication in scRNA-seq

Presenter: **Tuhin Majumder**, Wake Forest University, United States

Cell-cell communication (CCC) is often influenced by spatial proximity, prompting recent methods to incorporate spatial information using spatial transcriptomics (ST) data. However, in single-cell RNA-seq data, spatial coordinates are missing, limiting CCC inference to ligand-receptor expression patterns alone, potentially leading to incomplete or inaccurate conclusions. InterSpatial is introduced, a novel framework that leverages publicly available low-resolution ST data to enhance CCC analysis in single-cell data. InterSpatial performs cellular mapping, meta-cell clustering, and optimal transport to estimate inter-meta-cell distances. It examines ligand expression in sender meta-cells and trends of receptor expression in receiver meta-cells with respect to spatial distance, offering a visualization tool to screen for plausible CCC. Additionally, InterSpatial enables the application of ST-based CCC methods as downstream analyses. Its effectiveness is demonstrated by analyzing CCC between senescence-susceptible cells and macrophages in idiopathic pulmonary fibrosis (IPF) lungs, microglia and neurons in Alzheimer's disease brain tissue, and POSTN-positive fibroblasts and myeloid cells in myocardial infarction (MI) heart tissue. Results highlight the ability of InterSpatial to recover spatial context and uncover novel CCC patterns.

C0907: Multivariate principal component analysis for mixed-type functional data with application to mHealth in mood disorders

Presenter: **Debangana Dey**, Texas A&M University, United States

Mobile health studies collect various self-reported assessments capturing participants' behavior and well-being throughout the day. These assessments cover different scales, including continuous for physical activity, truncated for pain levels, ordinal for mood states, and binary for daily life events. Indexing these assessments by time and stacking them together, they form multivariate functional data with continuous, truncated, ordinal, and binary variables. A multivariate functional principal component analysis is proposed using a semiparametric Gaussian copula model, assuming a generalized latent non-paranormal process as the underlying mechanism. Latent temporal and inter-variable dependence is estimated through Kendall's Tau bridging method. The approach facilitates consistent function-on-function regression models, exemplified using data from 310 participants in the National Institute of Mental Health Family Study of the Mood Disorder Spectrum. The method characterizes the association between objectively collected physical activity and self-reported mood in individuals with major mood disorder subtypes, including Major Depressive Disorder and Type 1 and 2 Bipolar Disorder.

CO179 Room BCB 405 HIGH-DIMENSIONAL AND DYNAMIC BAYES: MODELS, COMPUTATION, APPLICATIONS Chair: David Rossell

C0356: Searching in ultra high dimensional sparse model spaces: New performance tests and boosting Gibbs sampling algorithms

Presenter: **Gonzalo Garcia-Donato**, Department of Economics and Finance - Universidad de Castilla La Mancha - Instituto de Desarrollo Regional, Spain

Co-authors: Maria Eugenia Castellanos Nueda

In the context of variable selection in sparse settings, a novel class of experiments is presented. These are based on the notion of contaminating real data sets with artificial spurious covariates in such a way that exact solutions can be easily computed. Such exact responses provide a direct and compelling way to evaluate the performance of search methods on model spaces of arbitrary cardinality. This tool is applied to Gaussian regression models, an important statistical problem that has benefited from the emergence of many new search methods in recent years. The contribution is also via revisiting classical Gibbs sampling algorithms, proposing new implementations that take advantage of sparsity. Despite their simplicity, the resulting methods are very competitive and fully automatic. A real genetic dataset is used to illustrate and motivate the various procedures presented in this research.

C0784: Bayesian temporal biclustering with applications to multi-subject neuroscience studies

Presenter: **Michele Guindani**, University of California Los Angeles, United States

Co-authors: Marina Vannucci, Erik Sudderth, Jaylen Lee, Megan Peters, Federica Zoe Ricci

The problem of analyzing multivariate time series collected on multiple subjects is considered, with the goal of identifying groups of subjects exhibiting similar trends in their recorded measurements over time as well as time-varying groups of associated measurements. To this end, a Bayesian model is proposed for temporal biclustering featuring nested partitions, where a time-invariant partition of subjects induces a time-varying partition of measurements. The approach allows for data-driven determination of the number of subjects and measurement clusters as well as estimation of the number and location of changepoints in measurement partitions. To efficiently perform model fitting and posterior estimation with Markov chain Monte Carlo, a blocked update of measurements' cluster-assignment sequences is derived. The performance of the model is illustrated in two applications to functional magnetic resonance imaging data and to an electroencephalogram dataset. The results indicate that the proposed model can combine information from potentially many subjects to discover a set of interpretable, dynamic patterns. Experiments on simulated data compare the estimation performance of the proposed model against ground-truth values and other statistical methods, showing that it performs well at identifying ground-truth subject and measurement clusters even when no subject or time dependence is present.

C0895: Birth-death dynamic graphical models

Presenter: **Elena Bortolato**, Universitat Pompeu Fabra, Spain

Co-authors: David Rossell, Stephen Hansen

In many real-world settings, from economics and biology to the social sciences, datasets often involve units that enter or disappear over time, leading to changing dimensionality and evolving dependence structures. Traditional graphical models struggle in these dynamic, potentially high-dimensional environments, especially when missingness arises not at random but from systematic birth death processes that alter the very compo-

sition of the observed system. We propose a Bayesian graphical model framework that explicitly incorporates the appearance and disappearance of units as a core part of the data-generating process, treating these structural changes as primary drivers of evolving dependencies. The model builds on recent advances in Bayesian factor and graphical modeling, employs a low-dimensional decomposition of the time-varying precision matrix, and incorporates latent states to capture distinct regimes, such as periods of stability or sudden structural shocks, through hidden Markov chains. This allows the framework to detect when and how dependencies among units fundamentally change. The result is a flexible tool for analyzing high-dimensional data with systematic structural missingness, providing robust inference on evolving dependence patterns when both the nodes and edges of the underlying graph change over time.

C0337 Room BCB 406 NOVEL METHODS FOR INFERRING NETWORK STRUCTURE
Chair: Keith Levin
C0859: Kernel methods for estimating distributions of graph statistics

Presenter: **Jonathan Stewart**, Florida State University, United States

Co-authors: Guang Qiu

A novel framework is introduced for estimating distributions of discrete network statistics using kernel density estimation methods. While empirical distributions provide valuable insights into network structures, they often suffer from sample-specific artifacts and irregularities that obscure relevant features of the true underlying distribution. The methodology bridges an important methodological gap in network science by extending discrete kernel methods to network domains without imposing strong parametric assumptions, and can provide more accurate estimates of distributions of graph statistics. The methodology applies to a diverse range of network characteristics, including degree distributions, shared partner counts quantifying triadic closure in a network, geodesic distances, and more. Theoretical properties of the smoothed estimators are established, and principled methods are introduced for bandwidth parameter selection through cross-validation that account for inherent dependence among edges in the network. Through simulation studies across different network types and sizes, the approach is demonstrated to substantially reduce estimation error, as measured by total variation distance, when compared with the empirical distribution, providing more reliable inference on the true underlying distribution. Lastly, the methodology is demonstrated in a network data application.

C0935: Asymptotically perfect seeded graph matching without edge correlation (and applications to inference)

Presenter: **Vince Lyzinski**, University of Maryland, College Park, United States

Co-authors: Tong Qi, Peter Viechnicki, Vera Andersson

The OmniMatch algorithm is presented for seeded multiple graph matching. In the setting of d -dimensional random dot product graphs (RDPG), it is proven that under mild assumptions, OmniMatch with s seeds asymptotically and efficiently perfectly aligns $O(s^x)$ unseeded vertices, for $x < 2d/4$, across multiple networks even in the presence of no edge correlation. The effectiveness of the algorithm is demonstrated across numerous simulations and in the context of shuffled graph hypothesis testing. In the shuffled testing setting, testing power is lost due to the misalignment/shuffling of vertices across graphs, and the capacity of OmniMatch is demonstrated to correct for misaligned vertices prior to testing and hence recover the lost testing power. The algorithm is further demonstrated on a pair of data examples from connectomics and machine translation.

C1069: Community detection via curvature gaps

Presenter: **Chang Li**, University of Virginia, United States

Co-authors: Zachary Lubberts, Melanie Weber, Yu Tian

The clustering problem is considered in stochastic blockmodel graphs from the perspective of Ollivier's Ricci Curvature, an extension of Ricci Curvature on manifolds to this discrete setting. The gap between the distributions of edge curvatures for within-cluster edges and between-cluster edges allows identifying these two groups of edges by their curvature, guaranteeing effective clustering. This curvature gap is studied under multiple signal strength regimes, identifying its limiting distribution and exploring the limits of curvature-based clustering. These distributional limits for edge curvatures are the first of their kind in the literature, and show that curvature-based clustering can be an effective competitor to traditional clustering methods, even in low signal strength settings.

C0166 Room BCB 407 BAYESIAN NONPARAMETRIC METHODS FOR LARGE-SCALE INFERENCE PROBLEMS
Chair: Linxi Liu
C1356: Bayesian nonparametric approaches for functional clustering

Presenter: **Wenyu Gao**, University of North Carolina at Charlotte, United States

Functional clustering of high-dimensional data is essential for uncovering latent structure in complex signals such as fMRI, where the number of underlying clusters is often unknown. Traditional functional clustering approaches, such as centroid-based, density-based, and parametric model-based methods, provide straightforward implementations but typically require pre-specifying the number of clusters and offer limited flexibility in capturing complex dependencies. In contrast, Bayesian nonparametric (BNP) methods, particularly Dirichlet process mixture (DPM) models, allow the number of clusters to be inferred directly from the data while naturally quantifying uncertainty in cluster assignments. Building on this framework, weighted Dirichlet process mixture (WDPM) models incorporate auxiliary or subject-level information through weights, enabling more informed clustering of functional signals. WDPM models are especially well-suited to high-dimensional neuroimaging studies such as fMRI in autism spectrum disorder (ASD), as they can accommodate high dimensionality, spatial correlation, and heterogeneity across subjects. This BNP framework offers a principled, data-driven alternative to conventional methods, providing enhanced flexibility, interpretability, and adaptability for functional clustering across diverse application domains.

C1351: Scalable multi-trait fine-mapping for metabolite GWAS

Presenter: **Weiqiong Huang**, Department of Statistics, University of Pittsburgh, PA, US, United States

Co-authors: Christopher McKennan, Joshua Cape, Emily Hector

Genome-wide association studies (GWAS) have enabled the discovery of thousands of genetic loci associated with complex traits, yet identifying the causal variants and understanding their biological mechanisms requires statistical fine-mapping. Existing fine-mapping approaches face two major limitations: They are often computationally prohibitive when scaling to hundreds of traits, and they typically assume independence of genetic effects, which can lead to biased inference in the presence of pleiotropy and linkage disequilibrium. A Bayesian factor model is introduced for scalable and interpretable multi-trait fine-mapping from GWAS summary statistics. The model decomposes genetic effects into shared components mediated by latent biological processes and trait-specific effects, allowing for the capture of both pleiotropic and trait-unique signals. Inference procedure leverages sequential approximation and efficient model space exploration, achieving orders-of-magnitude speedups over conventional methods while retaining statistical rigor. Crucially, the framework naturally accommodates biologically informed priors such as metabolic pathway structures in metabolomics to interpretable inference grounded in domain knowledge. Applying the method to metabolite GWAS of over 700 traits, findings from studies with 20 larger sample sizes are replicated, and novel causal variants that are inaccessible to current fine-mapping pipelines are uncovered.

C1526: A partial-likelihood approach to tree-based density modeling and its applications to Bayesian inference

Presenter: **Benedetta Bruni**, Duke University, United States

Co-authors: Li Ma

Tree based priors for probability distributions are specified using a predetermined, data independent collection of candidate recursive partitions of the sample space. To characterize a target density in detail, candidate partitions must expand deeply into all areas of the sample space with potential

non zero sampling probability. Such a system of partitions can incur prohibitive computational costs and cause overfitting, especially in regions with little probability mass. Existing models typically rely on relatively shallow trees. This hampers one of the most desirable features of trees, their ability to characterize local features, and reduces statistical efficiency. Traditional wisdom holds that this compromise is necessary for coherent likelihood based Bayesian inference, as a data dependent partition system that allows deeper expansion only in regions with more observations would induce double dipping of the data. We propose to restore coherency while allowing the candidate partitions to be data dependent, using Coxs partial likelihood. Our partial likelihood approach is applicable to existing likelihood based methods and to Bayesian inference on tree based models. We give examples in density estimation where the partial likelihood is endowed with existing priors on tree based models and compare with the standard, full likelihood approach. The results show substantial gains in estimation accuracy and computational efficiency from adopting the partial likelihood.

CO143 Room BCB 409 STATISTICS IN NEUROSCIENCE II
Chair: Julia Wrobel
C0754: Effect sizes variability in harmonization studies of Alzheimer's disease biomarkers

Presenter: **Dana Tudorascu**, University of Pittsburgh, United States

In recent years, multisite neuroimaging and blood biomarkers studies of Alzheimer's Disease have been rising in popularity due to increased statistical power and enabling the generalization of research outcomes; however, different blood biomarkers collection as well as summary neuroimaging measures derived from Positron Emission Tomography (PET) studies using different tracers and scanners, hinders the direct comparability of these measures across different sites. Furthermore, differences between these measures can lead to high variability in estimating effect sizes across clinical groups, depending on the blood assay, imaging scanner, or PET imaging tracer used. The aim is to present several standardization methods for blood biomarkers and for PET imaging-derived measures across different sites or scanners and to show effect-size variability in clinical groups using different harmonization techniques.

C1037: Hierarchical modeling of localized intracranial volume abnormalities in craniosynostosis

Presenter: **Joshua Lukemire**, Emory University, United States

Co-authors: Ryan Taylor, Julia Wrobel

Craniosynostosis, the premature fusion of cranial sutures, alters skull morphology and can restrict localized brain development despite normal overall intracranial volume. These localized deformations are linked to long-term neurodevelopmental deficits, but there is a current lack of quantitative models that capture the heterogeneity in growth patterns driven by different suture fusion types. Leveraging a large dataset of CT and 3D photogrammetry scans from craniosynostosis patients and healthy controls, a novel hierarchical generalized additive model that uses soap-film smooths is proposed, and basis functions are shared to model continuous volume maps across space, age, and fusion type. This approach enables fine-scale localization of volume abnormalities and prediction of future development trajectories. The model accounts for covariates such as age and sex and is estimated using a Bayesian framework.

C1040: Improving replicability of brain-behavior association studies by leveraging study design features

Presenter: **Kaidi Kang**, Wake Forest University School of Medicine, United States

Co-authors: Jakob Seidlitz, Richard Bethlehem, Jiangmei Xiong, Megan Jones, Arielle Keller, Ran Tao, Anita Randolph, Bart Larsen, Brenden Tervo-Clemmens, Oscar Miranda Dominguez, Jonathan Schildcrout, Damien Fair, Theodore Satterthwaite, A Alexander-Bloch, Simon Vandekar. Several recent studies raised concerns about the low replicability of brain-behavior association studies and showed that thousands of study participants are required for good replicability. However, massive sample sizes are often infeasible in practice. Analyses and meta-analyses are performed using 63 longitudinal and cross-sectional magnetic resonance imaging (MRI) studies from the Lifespan Brain Chart Consortium (77,965 total scans) to systematically investigate how the modifiable study design features can be leveraged to improve the ESs (and therefore, the replicability) of brain-behavior association studies. Based on strong empirical evidence and pragmatic statistical theory, concrete and actionable study design and analysis procedures are provided to neuroscientists to help them improve the replicability of their studies, given their different research objectives and the nature of their target associations.

CC435 Room BCB G09 MACROECONOMICS AND SOCIAL POLICY
Chair: Martin Wagner
C1362: The effects of interest rate increases on consumers' inflation expectations: The roles of informedness and compliance

Presenter: **Edward Knotek**, Federal Reserve Bank of Cleveland, United States

Co-authors: James Mitchell, Mathieu Pedemonte, Taylor Shiroff

How monetary policy communications associated with increasing the federal funds rate causally affect consumers' inflation expectations is studied in real time. In a large-scale, multi-wave randomized controlled trial (RCT), there is weak evidence that communicating these policy changes lowers consumers' medium-term inflation expectations on average. However, information differs systematically across demographic groups, in terms of ex ante informedness about monetary policy and ex post compliance with the information treatment. Monetary policy communications have a much stronger effect on the subset of consumers who had not previously heard news about monetary policy and who take sufficient time to read the treatment. The findings show that, in an inflationary environment, these consumers expect that raising interest rates will lower inflation. More generally, the results emphasize the importance of measuring both respondents' information sets and their compliance with treatment when using RCTs in empirical macroeconomics to better understand the real-world implications of monetary policy communications.

C0304: Escalated debt levels of Indian states: Who is responsible - the national or the state government?

Presenter: **Piyali Banerjee**, Ashoka University, India

The aim is to examine the co-movements and dynamics of state-level liabilities in India from 1994 to 2022 for 21 states. Using a dynamic factor model with time-varying loadings and stochastic volatility, the liabilities to GSDP ratio is decomposed into national, liability-specific, and state-specific factors. These factors help to assess how states with different liability levels and economic strength respond to fiscal pressures. The results reveal that common shocks, e.g., the GFC, Covid-19, and demonetization, play a dominant role in shaping the state liabilities. Over time, for the economically strong and low liability states, their influence has declined. Liability-specific factors, like the states with low and high liabilities, have gained prominence in the post-2014 era. The result is more evident in high GSDP states, indicating growing structural divergence. The empirical results show that states with strong fiscal discipline are characterized by higher own tax revenue and stable law and order. These states mostly rely on their internal debt levels. Besides, they are better insulated from national volatilities. In contrast, states heavily reliant on central transfers and subsidy-driven expenditure remain more exposed and are fiscally vulnerable. Findings highlight the need for a differentiated fiscal approach across states of India, with a policy focus on enhancing internal revenue capacity, rebalancing expenditure priorities, and fostering long-term fiscal autonomy.

C1378: Multidimensional panel data regression model: The case of the multidimensional homeownership vacancy rate in the USA

Presenter: **Talha Omer**, Linnaeus University, Vaxjo, Sweden

Co-authors: Daniel Henderson, Andros Kourtellis

The purpose is to extend the two-dimensional panel data regression model to a multidimensional setting for mixed-frequency data. Structured machine learning regression is explored in this context, incorporating techniques such as sparse group LASSO (sg-LASSO), LASSO, U-MIDAS, elastic net MIDAS, and average MIDAS. The theoretical properties and mathematical equations of these multidimensional panel data regression models are detailed. Their performance is tested via Monte Carlo simulations and two distinct data-generating processes (DGPs), each varying in

sample size, number of variables, and data frequency. As an empirical application, three-dimensional home ownership vacancy rates in the U.S. are nowcasted via the extended model, and the nowcasting robustness is assessed via the Diebold and Marino test. Additionally, a random walk forecast is computed and considered a benchmark for nowcast evaluation. The performance of all the MIDAS models used against traditional two- and three-dimensional panel data regression methods is compared using the out-of-sample root mean square error (RMSE) as a nowcasting accuracy metric for both simulated and empirical data. The results indicate that the three-dimensional MIDAS panel data regression model outperforms the comparative models, demonstrating a lower out-of-sample RMSE than the other models do.

CC424 Room BCB 212 CAUSAL INFERENCE AND POLICY EVALUATION
Chair: Ralf Wilke
C0212: Harnessing genetic variants for local average treatment effect estimation

Presenter: **Michela Bia**, Luxembourg Institute of Socio-Economic Research, Luxembourg

The contribution to the literature on genetic epidemiology and health economics is by estimating local average treatment effects (LATE) using genetic data. Using the Understanding Society dataset, a flexible instrumental variables (IV) framework is applied that leverages genetic variants to estimate heterogeneous effects of arthritis on equivalized income. Results show significant negative effects within some genetic subgroups and positive effects in others, possibly due to reliance on welfare programs among low-income individuals. This heterogeneity highlights the importance of subgroup analysis, which standard 2SLS may overlook. To improve the credibility and precision of the estimates, the testing approach was adopted by a recent study, which identifies valid instruments under the assumption that a relative majority is valid, conditional on compliance. This method is implemented using genetic data. Additionally, the findings are benchmarked with Mendelian randomization (MR), an alternative IV approach that supports the robustness of the results. New insights are offered into the economic effects of chronic conditions like arthritis, with relevant implications for public health and economic policy.

C0616: Optimal policy learning in empirical practice

Presenter: **Hannah Busshoff**, University of St. Gallen, Switzerland

Co-authors: Michael Lechner

The purpose is to review and extend recent methods in policy learning for heterogeneous treatment effects. Interpretable decision rules (decision trees) are compared with black-box approaches (policy forests and policy neural networks), highlighting trade-offs in model explainability, robustness to misspecification, stability, and computational demands. Using administrative data on active labor market policies, it is demonstrated how these methods can guide individualized treatment assignment. Results provide practical insights into the design of data-driven, transparent, and performant policy rules.

C0805: Resource-efficient policy targeting under heterogeneous partial interference

Presenter: **Laura Forastiere**, Yale University, United States

Co-authors: Elena Dal Torriente, Chan Park

In many empirical studies, units are interconnected, and a unit's outcome may depend on the treatment of others, leading to interference. When interference is heterogeneous, treating individuals with specific characteristics can influence the population average outcome differently, either through their direct response or their impact on others. For instance, policymakers may minimize resource use by vaccinating individuals identified as superspreaders to achieve a target reduction in disease incidence. Under heterogeneous clustered interference, a method to estimate optimal stochastic treatment allocations is proposed, in which an individual's treatment probability is determined by both individual- and cluster-level covariates. The approach minimizes the expected marginal treatment probability within a cluster while ensuring a specified outcome level is met. Although the resulting optimization problem is non-convex, it is efficiently solved using difference of convex functions algorithms. To evaluate the methodology, theoretical guarantees are provided, analyzing how the excess risk bound depends on the function class complexity and cluster size. Additionally, a simulation study is conducted, and the method is applied to a water, sanitation, and hygiene (WASH) intervention in Senegal. The estimated policy is compared to alternative approaches, with the method achieving greater resource efficiency compared to policies with homogeneous treatment probabilities within clusters.

CC396 Room BCB M201 STOCHASTIC PROCESSES
Chair: Nilanjan Chakraborty
C1395: A Gaussian-based approximation for a two-stage dual sourcing with delivery risk

Presenter: **Edson Antonio Goncalves de Souza**, IESEG School of Management/University of Lille, France

Co-authors: Stefano Nasini, Tanja Mlinar, Maud Van Den Broeke

The purpose is to develop a closed-form Gaussian approximation for a dynamic decision problem in inventory management with unreliable supply sources. The central contribution lies in establishing the analytical tractability of the value function through Gaussian partial moments and their integration into the optimality conditions. Specifically, closed-form expressions are provided for the expectations of the holding and shortage cost functions under normally distributed demand, enabling the bypass of simulation-based methods. Numerical experiments compare the Gaussian-based solution with policies derived from exponentially distributed demand. Across a broad set of parameterizations, the Gaussian approximation exhibits striking robustness: The resulting policies nearly coincide with those obtained under exponential assumptions, with average deviations negligible and maximum deviations diminishing as variability decreases. Importantly, the Gaussian formulation yields closed-form ordering rules that preserve the qualitative structure of the optimal policy while ensuring analytical simplicity. These results demonstrate that our methodology provides a powerful tool for decision-making under uncertainty, combining the interpretability of closed-form policies with accuracy comparable to simulation-driven methods, thereby offering a practically viable and computationally efficient framework for inventory decisions in settings with stochastic demand.

C1205: Price approximations and random sums

Presenter: **Jean Vaillancourt**, HEC Montreal, Canada

Co-authors: Bruno Remillard

Stock price approximation models, and more generally random sums, are used to infer knowledge about processes with discontinuous limits. Many difficulties arise when using the well-known metrics introduced by Skorohod on the space D of real-valued cadlag trajectories (right continuous with left limits everywhere). The space D is a natural choice for describing phenomena exhibiting bounded jumps, hence devoid of essential singularities (and no jump at the origin). Unfortunately, many results about the limiting behavior of such renewal counting processes in the literature are erroneous. An alternative metric resolves the issue, together with new central limit theorems for sequences of D -valued martingales.

C1197: A Riemannian covariance for manifold-valued data

Presenter: **Meshal Abuqrais**, King's College London, United Kingdom

Co-authors: Davide Pigoli

The extension of bivariate measures of dependence to non-Euclidean spaces is a challenging problem. The non-linear nature of these spaces makes the generalization of classical measures of linear dependence (such as the covariance) not trivial. The aim is to propose a novel approach to measure stochastic dependence between two random variables taking values in a Riemannian manifold, with the aim of both generalizing the classical concepts of covariance and correlation and building a connection to Frechet moments of random variables on manifolds. The purpose is to introduce generalized local measures of covariance and correlation, and it is shown that the latter is a natural extension of Pearson correlation. Then, suitable estimators are proposed for these quantities, and strong consistency results are proven. Finally, their effectiveness is demonstrated

through simulated examples and a real-world application.

CC440 Room BCB 408 APPLIED STATISTICAL ANALYSIS IN ECONOMICS AND SOCIETY
Chair: Tiejun Tong
C1382: Granular insights into food poverty: SAE-based estimates using HBS and healthy diet cost data

Presenter: **Stefano Marchetti**, Dipartimento di Economia e Management, Università di Pisa, Italy

Co-authors: Ilaria Benedetti, Haoran Yang

Ensuring the affordability of healthy and sustainable diets (HSDs) is a key challenge for food policy and social equity. The aim is to present a statistical framework to estimate the prevalence of food poverty across Italy's 107 provinces, using survey microdata and web-based price information. Household-level food consumption data from the Italian Household Budget Survey (HBS) 2022 are combined with monthly food basket costs derived from webscraped price data collected by the Ministry of Enterprises and Made in Italy. The food poverty indicator is defined by comparing household food expenditure to the theoretical cost of adhering to HSDs, adjusted by demographic needs. Due to limited sample sizes at the subnational level, a small area estimation (SAE) area-level approach is implemented to produce reliable provincial estimates. All computations were performed in the ADELE secure environment for official microdata. By means of spatially disaggregated policy indicators, actionable insights are provided for monitoring food affordability inequalities and guiding targeted interventions.

C0520: A hybrid smoothing approach for estimating airline passenger demand

Presenter: **Maria Rosa Nieto Delfin**, Investigaciones y Estudios Superiores, S.C, Mexico

Co-authors: Rafael Bernardo Carmona Benitez

The focus is on a hybrid model for estimating air passenger (pax) time series at the route level, thus addressing the limitations of existing unobserved components models, which primarily estimate trends. Accurate modeling of passenger demand is crucial for the air transportation industry, where time series exhibit strong seasonal behavior influenced by socioeconomic conditions and exceptional events. A methodology is proposed that integrates the trend-based unobserved components model developed by a prior study with the Holt-Winters multiplicative model. This approach enables the estimation of trend, seasonality, and cycle components in a unified framework. The proposed model involves the multiplication of the estimated trend component by the seasonal and cyclical components. This process enables the reconstruction of the original time series with greater precision. The hybrid model is validated using U.S. Bureau of Transportation Statistics route-level data, and its estimating accuracy is compared against the classical decomposition method and the Holt-Winters model. The performance of the model is assessed using the Diebold-Mariano test and mean absolute percentage error (MAPE). The findings indicate that the integration of a hybrid model enhances estimation accuracy by capturing the multi-component nature of pax time series. This approach provides a flexible and interpretable framework for time series in sectors with strong seasonal dynamics.

C1383: When cheap gets costly: Tracking cheapflation in Italy during high and moderate inflation

Presenter: **Luca Secondi**, University of Tuscia, Italy

Co-authors: Tiziana Laureti, Niccolo Salvini

The phenomenon of cheapflation - where the prices of low-cost goods rise faster than those of premium alternatives - has drawn increasing attention as a driver of inflation inequality. While it has been widely observed during recent episodes of sharp price acceleration, less is known about its persistence during more moderate phases of inflation. The aim is to examine cheapflation in Italy across both high and moderate inflation contexts, drawing on millions of web-scraped food prices combined with official household expenditure statistics. By exploring territorial patterns in the relationship between rising prices of budget goods and the concentration of vulnerable households, the regressive nature of inflation beyond aggregate indicators is investigated. The contribution is to show how the unequal burden of inflation can remain significant even as overall inflationary pressures ease, with important implications for monitoring and policy.

CV366 Room BCB G08 BIostatistics (VIRTUAL)
Chair: Alessia Pini
C0286: Structural equation modelling for diagnostic test accuracy meta-analysis

Presenter: **Zelalem Negeri**, University of Waterloo, Canada

Diagnostic test accuracy meta-analysis is a rapidly growing and active area of research. Standard approaches to this type of meta-analysis utilize hierarchical or bivariate random effects models that account for the within- and between-study heterogeneity in test characteristics. However, these approaches usually fail to converge for sparse data types or lead to biased inferences for meta-analysis with few studies. Moreover, the established methods cannot handle complex relationships between primary studies (e.g., direct and indirect effects or moderators and mediators), missing data, and adjusting for potential confounding variables. Therefore, the aim is to develop a structural equation modeling-based framework for aggregate data meta-analysis of diagnostic test accuracy studies to overcome these limitations of the standard methods. The proposed method is demonstrated and validated via extensive simulation studies and real-life data examples.

C1071: A user-friendly EEGLAB plug-in integrating functional data analysis to advance EEG research

Presenter: **Mohammad Fayaz**, Shahed University, Iran

EEG datasets pose complex challenges requiring advanced statistical methods. Functional data analysis (FDA) and machine learning offer powerful frameworks to model EEG signals over continuous domains. EEGLAB is a leading open-source EEG platform, widely cited in Web of Science and SCOPUS, and used globally, confirmed by bibliometric analyses. Despite FDAs potential, programming complexity limits its use among many researchers. The EEGLAB-FDA plug-in addresses this by integrating FDA methods into EEGLAB with a user-friendly graphical interface, enabling computational neuroscientists and statisticians to apply advanced FDA tools without coding. It includes modules for functional principal component analysis (FPCA), functional canonical correlation analysis (FCCA), and event-related potential (ERP) analysis, including smoothing, derivative estimation, and phase-plane visualization. Supporting dense and sparse data, recent updates add functional regression and enhanced sparse data handling. Bug fixes and interface improvements improve usability. Applications with visual and auditory oddball tasks show the plug-in's capacity to reveal insights, especially when analyzing the P300 component. Available free on GitHub and the EEGLAB plug-in list, it is supported by a YouTube channel with tutorials. Ongoing development aims to expand its features and impact in computational neuroscience and statistics.

C1241: Statistical testing in longitudinal studies for diffusion tensor imaging

Presenter: **Lyudmila Sakhanenko**, Michigan State University, United States

Co-authors: Juna Goo, David Zhu

A longitudinal diffusion tensor imaging (DTI) study on a single brain can be remarkably useful to probe white matter fiber connectivity that may or may not be stable over time. The aim is to consider a novel testing problem where the null hypothesis states that the trajectories of a coherently oriented fiber population remain the same over a fixed period of time. Compared to other applications that use changes in DTI scalar metrics over time, this test focuses on the partial derivative of the continuous ensemble of fiber trajectories with respect to time. The test statistic is shown to have the limiting chi-square distribution under the null hypothesis. The power of the test is demonstrated using Monte Carlo simulations based on both the theoretical and empirical critical values. The proposed method is applied to a longitudinal DTI study of a normal brain.

Monday 15.12.2025

10:30 - 12:10

Parallel Session L – CFE-CMStatistics 2025

CI007 Room BCB 206 CFE SPECIAL INVITED SESSION: A TRIBUTE TO H. PESARAN II**Chair: Andrew Harvey****C0166: M. Hashem Pesaran: Some reminiscences and reflections***Presenter:* **Richard J Smith**, University of Cambridge, United Kingdom

Since our initial acquaintance in around the mid-1980s Hashem and I have been co-authors and remain colleagues and friends. I will first briefly review our research collaborations over the years through joint work from the mid 1980s to the early 2000s; viz. "A Unified Approach to Estimation and Orthogonality Tests in Linear Single Equation Econometric Models", *Journal of Econometrics*, 44 (1990), 41-66; "A Generalized R2 Criterion for Regression Models Estimated by Instrumental Variables", *Econometrica*, 62 (1994), 705-710; "Structural Analysis of Vector Error Correction Models with Exogenous I(1) Variables", *Journal of Econometrics*, 97 (2000), 293-343; "Bounds Testing Approaches to the Analysis of Long-Run Relationships", *Journal of Applied Econometrics*, 16 (2001), 289-326. I will then discuss his importance to and influence on the Faculty of Economics and Politics at the University of Cambridge during the first period 1989-95 when we first overlapped and then since 2005 when I re-joined what had become the Faculty of Economics.

C0161: Estimation and inference in high-dimensional panel data models with interactive fixed effects*Presenter:* **Oliver Linton**, University of Cambridge, United Kingdom*Co-authors:* Maximilian Ruecker, Michael Vogt, Christopher Walsh

New econometric methods are developed for estimation and inference in high-dimensional panel data models with interactive fixed effects. The approach can be regarded as a non-trivial extension of the very popular common correlated effects (CCE) approach. A projection device is first constructed to eliminate the unobserved factors from the model by applying a dimensionality reduction transform to the matrix of cross-sectionally averaged covariates. The unknown parameters are then estimated by applying lasso techniques to the projected model. For inference purposes, a desparsified version of the lasso-type estimator is derived. While the original CCE approach is restricted to the low-dimensional case where the number of regressors is small and fixed, methods can deal with both low- and high-dimensional situations where the number of regressors is large and may even exceed the overall sample size. Theory for the estimation and inference methods is derived both in the large- T -case, where the time series length T tends to infinity, and in the small- T -case, where T is a fixed natural number. Specifically, the convergence rate of the estimator is derived, and it is shown that its desparsified version is asymptotically normal under suitable regularity conditions. The theoretical analysis is complemented by a simulation study and an empirical application to characteristic-based asset pricing.

C0155: Hidden threshold models with applications to asymmetric cycles*Presenter:* **Andrew Harvey**, University of Cambridge, United Kingdom*Co-authors:* Jerome Simons

Threshold models are set up so that there is a switch between regimes for the parameters of an unobserved components model. When Gaussianity is assumed, the model is handled by the Kalman filter. The switching depends on a component crossing a boundary, and, because the component is not observed directly, the error in its estimation leads naturally to a smooth transition mechanism. A prominent example motivating thresholds is that of a cyclical time series characterized by a downturn that is more or less rapid than the upturn. The situation is illustrated by fitting a model with three potentially asymmetric cycles, each with its own threshold, to observations on ice volume in Antarctica since 799,000 BCE. The model is able to produce multi-step forecasts with associated prediction intervals. A second example shows how a hidden threshold model is able to deal with the asymmetric cycle in monthly US unemployment.

CO140 Room BCB G07 HiTEC: RECENT ADVANCES IN MODEL SPECIFICATION TESTING**Chair: Bojana Milosevic****C1097: Closure of noncentral Wishart mixtures and testing random effects in multivariate factorial designs***Presenter:* **Frederic Ouimet**, Université du Québec à Trois-Rivières, Canada

Consider mixtures of noncentral Wishart distributions where both the mixing and mixed laws share the same degrees of freedom. It is shown that such a mixture remains noncentral Wishart with the same degrees of freedom, a closure property that extends known univariate chi-square results to arbitrary dimension $d \geq 1$ and general scale matrices. Leveraging this fact, exact finite-sample reference distributions are derived for multivariate tests of random effects in balanced two-factor factorial designs with d -dimensional normal responses. In particular, the classical eigenvalue-based statistics built from sums of outer products and the pooled error matrix have a matrix-variate beta type II (matrix F) distribution under the null, enabling rigorous tests for covariance components without asymptotic approximations. The methodology applies as well when some factors are fixed, and extends readily to models with more than two factors.

C1102: Tensor changepoint detection and eigenbootstrap*Presenter:* **Michal Pesta**, Charles University, Czech Republic*Co-authors:* Barbora Pestova, Martin Romanak

Tensor data consisting of multivariate outcomes over the items and across the subjects with longitudinal and cross-sectional dependence are considered. A completely distribution-free and tweaking-parameter-free detection procedure for changepoints at different locations is designed, which does not require training data. A CUSUM-type test statistic is employed, and its asymptotic properties are derived for a large number of available individual profiles. The considered test is shown to be consistent. The aim is to propose eigenbootstrap superstructure that overcomes the computational curse of dimensionality without any loss of information, while it preserves all the dependencies within and between the panels. The validity of this new and fast resampling algorithm is proved in this general setting. The empirical properties of the detection technique are investigated through a simulation study. The fully data-driven test is applied to real-world data from EEG and psychometrics.

C1103: On the number of replications in resampling tests and Monte Carlo simulation studies*Presenter:* **Daniel Gaigall**, FH Aachen University of Applied Sciences, Germany*Co-authors:* Julian Gerstenberg

The purpose is to investigate rejection probabilities of statistical tests based on resampling procedures. The general framework under consideration covers, in particular, bootstrap and permutation techniques. It turns out that specific properties of the P-value distribution play a key role, namely convexity or concavity, the Bernstein property, and those of beta mixture models. A detailed analysis is provided, and it is clarified how these properties relate to each other. New bounds are derived for the rejection probability. The results link the number of replications with size and power of the test. Numerical considerations demonstrate the quality of the bounds. An important application is the nested simulation estimator in Monte Carlo simulation studies. Findings indicate that a moderate or even rather small number of replications is sufficient to obtain useful simulation results. This enables a substantial reduction of the computational effort in Monte Carlo simulation studies.

C1157: On the multi-sample problem in the presence of random right censoring*Presenter:* **Jaco Visagie**, North-West University, South Africa*Co-authors:* James Allison, Anke Steyn

A classical problem in survival analysis is testing whether two or more independent samples are realized from the same distribution. In practice, this problem is often complicated by the presence of random right censoring. A new test statistic is introduced based on empirical characteristic

functions. The proposed test statistic is a weighted L2-type distance between the empirical characteristic functions of the observed samples. The counterpart of the empirical characteristic function in the presence of censoring is based on the Kaplan-Meier estimator. The finite sample properties of the newly proposed test are investigated via simulated as well as observed data.

CO165 Room BCB G08 ADVANCED STATISTICAL METHODS FOR ENERGY ECONOMICS
Chair: Giuseppe Scandurra
C0644: Measuring fuzzy energy poverty in Italy

Presenter: **Giuseppe Scandurra**, Parthenope University of Naples, Italy

Energy poverty remains a critical yet inadequately captured dimension of socio-economic deprivation, especially in high-income countries, where affordability, rather than access, defines the problem. Traditional measurement methods offer binary classifications that overlook the gradient and complexity of energy-related hardship. These shortcomings are addressed by proposing a novel, fuzzy logic-based approach to assessing energy poverty. Building on the totally fuzzy and relative framework, the aim is to propose the fuzzy energy poverty index (FEPI) that integrates both objective and subjective indicators, including energy expenditure-to-income ratios, thermal comfort adequacy, housing insulation quality, and self-reported hardship. Fuzzy membership functions are employed to quantify the degree of deprivation for each dimension, which are then aggregated using fuzzy operators to form a composite index. Applied to empirical data on Italian households, the FEPI captures a continuum of energy deprivation, highlighting not just those in acute need but also those at risk, offering more nuanced insights than binary classifications. This approach enables dynamic responsiveness to socio-economic and environmental pressures such as inflation, energy price volatility, and climate-related demands on residential energy use. FEPI model contributes to more targeted, inclusive, and sustainable energy policy interventions, aligning with European goals for energy justice and social equity.

C0739: Fiscal incentives in Italian residential sector: Long-term effectiveness and distributional effects

Presenter: **Alfonso Carfora**, University of Macerata, Italy

Fiscal incentives' effectiveness is examined for energy efficiency and renewable transition in Italy's residential sector through regional analysis. Using panel vector autoregressive (PVAR) models with panel-specific fixed effects on data from 20 Italian regions (2009-2021), the impact of dwelling renovation bonuses (VDR) and energy efficiency bonuses (VEE) on energy consumption patterns is investigated. Findings reveal fiscal incentives demonstrate gradual but significant energy consumption reduction, with effects materializing after several years and stabilizing around zero after 4-5 years. This suggests sustained long-term policy commitment is essential for meaningful efficiency improvements. While residential measures show promise for consumption reduction, their effectiveness in driving broader fossil-to-renewable transition is limited, indicating the need for comprehensive multi-sectoral approaches. Using non-hierarchical fuzzy clustering to explore regional heterogeneity, we reveal significant geographical disparities. Northern regions capture over 60% of deductions, highlighting distributional concerns and the regressive nature of current structures that disproportionately benefit high-income homeowners. Results demonstrate that while fiscal incentives like Superbonus 110% can drive efficiency improvements, design requires fundamental reform addressing equity concerns and enhancing effectiveness for vulnerable groups, including low-income households and tenants.

C0803: Energy efficiency and household energy poverty: The impact of rebound effects in Italy

Presenter: **Leo Fulvio Minervini**, University of Macerata, Italy

Post-pandemic energy price volatility has intensified affordability challenges, prompting Italian policymakers to implement initiatives like Superbonus 110% and to extend the existing Ecobonus to promote residential energy efficiency (EE). However, the effectiveness of these interventions in reducing energy poverty (EP) remains unclear, particularly when considering behavioral rebound effects (RE) that may offset expected benefits. The relationship is empirically examined between energy efficiency improvements and energy poverty reduction using micro-level data from the 2023 ISTAT Household Budget Survey. Methodologically, it employs a two-stage procedure combining a stochastic energy demand frontier model with a probabilistic approach to quantify household-level energy efficiency and direct rebound effects. Results reveal that while energy efficiency improvements significantly reduce energy expenditures and EP risk, these benefits are substantially eroded by behavioral rebound effects, especially among households with initially low efficiency levels. Findings provide crucial insights for designing more effective energy efficiency policies to combat energy poverty in economically stratified contexts.

C0192: Ensuring the security of the transition: Examining the impact of geopolitical risk on the price of critical minerals

Presenter: **Jamel Saadaoui**, University of Paris 8, France

Co-authors: Russell Smyth, Joaquin Vespignani

Constant and time-varying parameter local projection (TVP-LP) regression models are used to examine the effect of geopolitical risk on prices of six critical minerals: Aluminum, copper, nickel, platinum, tin, and zinc. A conceptual framework is proposed in which the responsiveness of prices for critical minerals to geopolitical risk depends on the non-technical risk associated with procuring each critical mineral, and geopolitical threats have a bigger effect on critical mineral prices than geopolitical acts. Results are generally consistent with these predictions. Considerable evidence is found that the effect of geopolitical risk on the prices of critical minerals is time-varying, with the Gulf War, 9/11 terrorist attacks, and COVID-19 pandemic each having a significant effect. It is found that shocks due to geopolitical threats are generally bigger in magnitude than geopolitical acts and that prices respond more quickly to geopolitical threats.

CO220 Room BCB G09 CLEVER MODELS FOR COMPLICATED DATA
Chair: Sondre Hoelleland
C0444: Extending landmarking to mixture cure models with longitudinal covariates

Presenter: **Marta Cipriani**, Sapienza University of Rome, Italy

Co-authors: Marco Alfo, Mirko Signorelli

Dynamic prediction models represent an essential class of models for personalized medicine, providing real-time updates on prognosis based on evolving patient information. Among these, the landmarking approach has gained popularity due to its flexibility and conceptual simplicity. However, its integration into cure models remains underexplored. In the context of mixture cure models (MCM), current applications of landmarking rely exclusively on traditional summarization techniques for time-varying covariates, notably based on the last observation carried forward approach. A novel dynamic prediction framework is proposed that extends model-based landmarking to MCMs. The framework separates prediction into two components: (1) the incidence component, which estimates the probability of being uncured using logistic regression and baseline covariates, and (2) the latency component, which estimates post-landmark survival among uncured individuals through a Cox proportional hazards model that incorporates summaries of longitudinal data trajectories. Specifically, the longitudinal trajectories of patient covariates are modeled up to the landmark time using linear mixed-effects models or multivariate generalized linear mixed-effects models. These models allow for the estimation of individual-specific random effects, which provide a compact and informative summary of the patient's covariate trajectory and are then used as (fixed) predictors in the cure model.

C0540: Sunset, the dwelling-specific condition, and socio-economic drivers of housing prices in a Norwegian urban area

Presenter: **Ingrid Sandvig Thorsen**, University of Bergen, Norway

Co-authors: Baard Stove, Kristian Gundersen

The effect of sunshine and the condition of real estate in the Norwegian housing market is studied. The correlation between local housing prices and the socioeconomic makeup of a neighborhood is also studied. Spatial Gaussian Markov random field (GMRF) models are used to account for dependence between observations, also incorporating a set of other attributes as standard covariates. The sunshine property experience is based

solely on the landscape aspect, i.e., buildings and obstacles that potentially cast shadow on the property are not considered. In the Norwegian housing market, a condition report of the real estate is mandatory for a legal transaction. The condition of a house is assessed by a professional appraiser according to regulations, and different parts of the building are classified into four categories, which in total describe the condition of the real estate. Using house prices from the Bergen municipality, Norway, 2016-2022, the well-known relations are found in hedonic price models, concerning, for instance, age and size of house, lot size, and socioeconomic variables. However, it is also found that the condition of a house affects its price; a positive sunshine effect is documented, and the relationship between local socioeconomic conditions and housing prices is quantified. This is done in a model accounting for the spatial random field. To adjust for spatial confounding issues, a restricted spatial regression model (RSR) is used.

C0646: Innovation in high-dimensional categorical Bayesian optimization

Presenter: **Timothee Bacri**, University of Exeter, United Kingdom

Co-authors: Daniel Williamson, Bertrand Nortier

Bayesian optimization is a powerful tool to optimize black-box and expensive objective functions. Handling high-dimensional inputs is challenging, but can be managed with methods such as embeddings. This allows the usage of traditional Gaussian processes as surrogates. A high-dimensional categorical optimization problem is looked into, and hence, discretization is required. However, the popular discretization method scaling-and-rounding suffers when dealing with large numbers of categories due to the non-uniform distribution of the input variables. An adaptation is proposed using quantile-binning instead. Bins being defined with quantiles means each category is equally likely over each variable, ensuring a balanced exploration over the high-dimensional categorical space. This method is evaluated, and its strengths and limitations are highlighted.

C0707: Robust hidden semi-Markov models for meteorological residuals in the Venice lagoon

Presenter: **Lorena Ricciotti**, University of Bari, Italy

Co-authors: Alfonso Russo, Sondre Hoelleland, Antonello Maruotti

The aim is to present a flexible statistical framework to analyze the tides' meteorological component in the Venice lagoon. A robust class of hidden semi-Markov regression models (HSMRMs) is proposed, capturing key marine data features such as regime switching, time-varying heteroskedasticity, heavy tails, skewness, and outliers. The framework extends hidden Markov regression models, relaxing the geometric sojourn time assumption and incorporating variable selection via elastic-net regularization. To improve robustness to outliers and skewed data, Gaussian, Student-t, and Johnson's SU distributions are used, enabling more accurate modeling of sea level behavior. Empirical analysis is conducted using hourly data from the Lido di Porto inlet tide gauge, covering meteorological variables as wind speed and direction, air and water temperature, pressure, and humidity. The proposed model identifies four environmental regimes influencing meteorological residuals associated with specific weather conditions and temporal dynamics. The model distinguishes between true and apparent contagion by incorporating autoregressive components, thus addressing latent regime shifts and explicit temporal dependencies. Simulation studies confirm the ability of the Johnsons SU-based HSMRM in parameter estimation accuracy and classification performance under skewed data generation processes. The regularization approach effectively selects relevant covariates and lags, enhancing model interpretability.

CO333 Room BCB 207 CONTEMPORARY ISSUES IN ECONOMICS AND FINANCE

Chair: Tanakorn Likitapiwat

C0301: Participation, transparency, and cost asymmetry in common-pool resource governance

Presenter: **Pathomwat Chantarasap**, Chulalongkorn University, Thailand

Co-authors: Nuttaporn Rochanahastin

The purpose is to explore the governance of common-pool resources, focusing on fisheries, where not all players/stakeholders are involved in setting usage rules. Motivated by real-world challenges in inclusive decision-making, a computer-based laboratory experiment is relied on to investigate how resource use is affected by different rule-setting scenarios. It is found that externally set resource-use rules (such as from government agencies) lead to resource usage that is closer to the social optimum than rules solely set by users. Additionally, the involvement of only a subset of players in rule setting significantly affects common resource use, particularly when cost structures differ among users. Players who are excluded from the decision-making process, especially those with lower extraction costs, tend to exploit the resource more aggressively, resulting in both overuse and increased inequality. These findings highlight the interaction between institutional exclusion and economic heterogeneity. By systematically testing these dynamics, robust empirical evidence is provided for the critical role of inclusive and transparent rule setting in promoting equitable participation and sustainable resource management.

C0315: Culture counts: Exploring firm performance through machine-learning based corporate culture in Thailand

Presenter: **Tanakorn Likitapiwat**, Chulalongkorn University, Thailand

A comprehensive Thai-specific corporate culture dictionary is constructed, using word embedding machine learning techniques, and the relationship is explored between the text-based corporate culture and firm performance for all Thai listed firms over the period of 2000 to 2021. A positive effect of corporate culture is found on firm performance such that a one standard deviation rise in corporate culture improves firm performance by 9%. Results also demonstrate that all five elements of corporate culture significantly drive Thai firms' performance; However, the effect is in the following order from highest to lowest: Teamwork, respect, innovation, quality, and integrity. Additional analyses also indicate that a strong culture, except for innovation and integrity, insulates Thai firms from the adverse effects of the 2008 global financial crisis on firm performance. Findings are robust to various tests, including a two-stage least squares instrumental variable approach, propensity score matching, entropy balancing, and lagged regression analysis. Finally, business cultures unique to the Thai context are highlighted, prioritizing a corporate culture centered on teamwork and respect, as opposed to Western culture, which places emphasis on innovation and integrity.

C1019: Saving and dissaving behavior in an aged society

Presenter: **Nuttaporn Rochanahastin**, Prince of Songkla University, Thailand

The purpose is to examine saving and dissaving behaviors across different age groups and generational cohorts in Thailand, using nearly three decades of repeated cross-sectional data from the Thai Household Socio-Economic Survey (HSES). The findings reveal a clear hump-shaped life-cycle pattern in saving behavior, with savings peaking between the ages of 56 and 65, slightly beyond traditional working life. Importantly, wealth accumulation remains positive even into later life stages, suggesting the influence of precautionary motives, cultural bequest norms, and limited annuitization options. Generational comparisons show that Baby Boomers and Generation X consistently save more than Generation Y, reflecting differences in economic experiences and structural opportunities across cohorts. The impact of economic experiences and life stages on saving behavior is underscored. These findings highlight the critical interplay of temporal, demographic, and cohort effects, offering valuable insights for policymakers seeking to promote financial resilience and security in an aging society.

C1096: Geopolitical and tourism risks: A comparison of analyses between econometric approaches and machine learning

Presenter: **Chanon Thongtai**, Chulalongkorn University, Thailand

Co-authors: Nuti Sornnil

The aim is to test predictive models and tools between econometric and machine learning methods, as well as to examine whether influence geopolitical risks influence the forecast of tourism demand in Taiwan. Monthly-based data is collected with a total of 264 observations with 8 variables from January 1998 to December 2019, before the COVID pandemic, to provide a shock-free dataset. Various models are used, including: WLS linear regression, k-nearest neighbors, decision tree, random forest, and gradient boosting with metrics as values and R^2 and root mean square

error (RMSE). The results show that the decision tree model is more effective in forecasting than other types of models. Decision tree is the model with the best indicator values: $R^2 = 0.42$ and RMSE at 0.0899. The econometric analysis suggests that geopolitical risks have a very small, yet negative, effect on the forecast of tourism demand.

CO324 Room BCB 208 INNOVATIONS IN BAYESIAN NETWORK PSYCHOMETRICS
Chair: Maarten Marsman
C0497: How much data is enough? Effective sample size in Bayesian graphical models

Presenter: **Giuseppe Arena**, University of Amsterdam, Netherlands

Co-authors: Lourens Waldorp, Maarten Marsman

In Bayesian analysis, the concept of effective sample size (ESS) applies separately to both the data and the prior, quantifying the independent information provided by the data beyond what is already encoded in an informative prior derived from a previous study. In graphical models such as Gaussian graphical models (GGMs), calculating the ESS becomes challenging due to the complexity of the model and the dependencies among parameters. While extensive prior work has introduced methods for estimating the ESS of data and prior, their use in the context of Bayesian graphical models remains relatively unexplored. The methodology for calculating the ESS of data and informative priors in GGMs is first outlined. Two common research scenarios are then considered: Planning a Bayesian analysis with a known informative prior and calibrating priors to prevent prior dominance in the posterior. In the former, the ESS quantifies the amount of data needed to balance prior influence, supporting robust inference. In the latter, the ESS is used to evaluate prior dominance, guiding the appropriate adjustments to maintain reliable posterior inference. The discussion concludes with methodological extensions for estimating the ESS in ordinal Markov random fields.

C0500: A stochastic block prior for clustering in graphical models

Presenter: **Nikola Sekulovski**, University of Amsterdam, Netherlands

Co-authors: Giuseppe Arena, Jonas Haslbeck, Karoline Huth, Nial Friel, Maarten Marsman

Existing statistical methods for analyzing graphical models in psychology often ignore the assumption of clustering, which refers to the grouping of variables that are more densely connected, despite its relevance in many psychological theories. The stochastic block model (SBM) is proposed as a prior distribution on the network structure in models for binary and ordinal data. The SBM assumes that variables belong to latent clusters and that the probability of an edge depends on cluster membership. Embedding the SBM in a Bayesian graphical modeling framework enables the formal incorporation of theoretical expectations about clustering, testing hypotheses about the number of clusters, and estimating cluster membership of the nodes from cross-sectional data. The benefits of this approach are demonstrated through a simulation study and a reanalysis of 30 empirical datasets. This method provides a principled approach to latent cluster inference in psychological network analysis by incorporating structural assumptions directly into the model through the prior.

C0806: A comparison of variable selection algorithms with an application to the ordinal Markov random field

Presenter: **Don van den Bergh**, University of Amsterdam, Netherlands

Bayesian variable selection plays a crucial role in network models, especially because the space of models is too large to enumerate. A range of Bayesian variable selection algorithms is compared in terms of their ability to accurately recover posterior inclusion probabilities and their speed of convergence. First, the performance of reversible jump MCMC, mixtures of mutually singular distributions, Rao-Blackwellized estimators, and sticky piecewise deterministic Markov processes is contrasted on a toy regression problem where enumeration is feasible. The methods are compared on accuracy and computational speed. Next, the comparison is extended to a more complex model where enumeration is infeasible: The ordinal Markov random field. In this network model, the number of parameters subject to selection grows quadratically with the number of variables, making it an excellent example to study how the methods scale.

C0814: Bayesian estimation of ordinal cross-lagged panel model

Presenter: **Vipasha Goyal**, University of Amsterdam, Netherlands

Co-authors: Maarten Marsman

The cross-lagged panel model (CLPM) is a widely used statistical method in psychological research for examining reciprocal and dynamic relationships between variables over time and drawing causal inferences. Despite their widespread use, existing implementations of the CLPM are limited in their capacity to support flexible Bayesian estimation and inference, restricting our ability to quantify model uncertainty. Without a formal treatment of model uncertainty, parameter estimates are conditioned on one selected model, which can lead to overconfident and biased inferences when multiple plausible models exist for the data at hand. Moreover, most implementations assume continuous and normally distributed variables, whereas many psychological constructs are measured using ordinal scales. To address these limitations, a Bayesian framework for ordinal CLPM is developed, incorporating Bayesian model averaging techniques for variable selection. This framework enables efficient estimation of dynamic relationships over time through posterior inclusion probabilities and inclusion Bayes factor. The proposed models and procedures are made available in free software packages in R and JASP.

CO319 Room BCB 210 CROSS-SECTIONAL ASSET PRICING
Chair: Paolo Zaffaroni
C0664: Low-frequency risk factors and their fundamental drivers

Presenter: **Sicong Li**, The Chinese University of Hong Kong, Hong Kong

There is a zoo of factors that capture systematic risk premia and a large number of economic variables that explain their time variation, which poses a doubly high-dimensional challenge to understanding how economic fundamentals relate to the time-varying dynamics of risk premia. A method is proposed to regularize this problem by identifying low-frequency risk factors, whose risk premia are driven by latent low-frequency state variables. Empirically, one below-business-cycle-frequency factor and one business-cycle-frequency factor, whose variation concentrates on cycles longer than eight years and between 1.5 and eight years, explain the expected returns of individual stocks and characteristic-managed portfolios. The below-business-cycle-frequency factor has a high Sharpe ratio, and stocks whose current size is small compared to their long-term average load on it. Moreover, selected macroeconomic and financial variables have statistically and economically significant out-of-sample predictive power for the returns of the two low-frequency factors.

C0993: Simple out-of-sample tests for asset pricing

Presenter: **Svetlana Bryzgalova**, London Business School, United Kingdom

Co-authors: Alberto Quaini, Ashish Sahay

The aim is to show that traditional measures of out-of-sample model performance in asset pricing ignore model estimation risk and significantly underscore true standard errors. As a result, typical tests overestimate t-stats associated with out-of-sample Sharpe ratios, alphas, etc. A simple split-sample estimation design is proposed, that allows to effectively measure out-of-sample model performance and provides valid statistical inference for both in-sample and out-of-sample parameters. Empirically, the performance of popular linear factor models is revisited, and it is found that model estimation risk has a nontrivial impact on the out-of-sample tests. The results have important implications for the evaluation of asset pricing models (linear, nonlinear, and those estimated via machine learning techniques), the use of spanning tests in model comparison, and measuring risk-adjusted returns out of sample.

C1109: International investing: Diversification and beyond

Presenter: **Soohun Kim**, KAIST, Korea, South

Co-authors: Andreas Neuhierl, Robert Korajczyk

The purpose is to develop a framework for analyzing individual stocks in a foreign market and constructing effective trading strategies for a home-country investor. The approach exploits (i) mispricing in the foreign market, (ii) risk premia disparities between home and foreign countries, and (iii) foreign country-specific investment opportunities. Applying this framework to 27 developed and developing countries, evidence of substantial international investment opportunities is provided, even among G7 countries. In particular, the profitability of strategies exploiting risk premia disparities challenges the prevailing view of a tightly integrated international financial market.

C1064: Market-based incentives for optimal audit quality

Presenter: **Huaizhi Chen**, University of Texas at Dallas, United States

The purpose is to examine how equity markets respond to the public release of audit-firm inspection reports by the U.S. regulator. Investors react differently based on the identifiability of the public issuers whose audits are covered in the inspection report. Auditors with identifiable issuer clients show positive abnormal returns for non-deficient reports and negative reactions for deficient ones. In contrast, issuers less easily linked to specific auditor inspections experience muted responses. More timely publication of inspection reports intensifies market reactions, while delays reduce their informativeness. The findings highlight how regulatory transparency can enable investors to better incorporate audit quality information into equity prices. Implications are discussed for market-based incentives for issuers and auditors.

CO072 Room BCB 211 TIME SERIES ECONOMETRICS

Chair: Antonio Montanes

C0649: Macroeconomic effects of temperature distributional shocks

Presenter: **Lola Gadea**, University of Zaragoza, Spain

Co-authors: Jesus Gonzalo, Andrey Ramos

The aim is to propose a novel methodology to assess the macroeconomic effects of distributional temperature shocks, persistent deviations in temperature quantiles from long-term trends, identified using the Hamilton filter. As these shocks are correlated across regions, a factor model is employed to extract their common components. Empirical analysis at the global level, as well as for the United States and the Euro Area, identifies three orthogonal factor shocks, each with significant macroeconomic implications. The first factor, which captures a general shift in the temperature distribution, leads to persistent declines in output and total factor productivity, consistent with prior studies focusing on average temperature. The main contribution is to uncover two additional factors that are uncorrelated with average temperature but affect specific parts of the distribution. These factors generate additional, previously undocumented macroeconomic effects. For instance, the third factor, which alters both tails of the distribution in the same direction, significantly reduces output and productivity in the Euro area. Findings underscore the importance of analyzing the full distribution of temperature shocks, offering new insights for estimating the social cost of carbon and assessing climate-related economic risks.

C0576: Modelling climate heterogeneity using an unconditional quantile vector error correction approach

Presenter: **Jesus Gonzalo**, Universidad Carlos III de Madrid, Spain

Co-authors: Lola Gadea, Andrey Ramos

Understanding the uneven spatial and temporal patterns of climate change is essential for assessing its broader societal and economic implications. A novel quantitative framework is presented based on an unconditional quantile vector error correction model (QVECM). In this approach, the unconditional quantiles of temperature and greenhouse gas concentrations are modeled as co-trending variables, consistent with the dynamics of a standard climate one-dimensional energy balance model (1-EBM). The QVECM framework enables: Estimation of diverse climate sensitivity parameters across the temperature distribution; Identification of warming amplification effects; Introduction of a Quantile Permanent-Transitory Decomposition (QPT), which separates the long-term drivers of climate change from short-term fluctuations. The permanent component identified by the QPT highlights the fundamental forces driving persistent warming trends. These insights are critical for informing the design of effective climate mitigation and adaptation strategies.

C0554: Robust estimation of the autocorrelation function in the presence of outliers.

Presenter: **Antonio Montanes**, University of Zaragoza, Spain

A robust and computationally efficient method is proposed for estimating the autocorrelation function (ACF) in time series data contaminated by outliers. While traditional estimators of the ACF perform well under normality, they are known to be highly sensitive to non-normality and outliers, leading to biased results. The Hurwicz estimator is revisited, which has been shown to be median unbiased under normal distributions and generalized error distributions. By interpreting the first-order autocorrelation as a ratio of centered normal variables, a robust estimator is derived based on the truncated Cauchy distribution, mitigating the issue of undefined moments inherent to the standard Cauchy. The methodology is extended to higher-order autocorrelations using pairwise transformations and robust statistics, maintaining computational simplicity. Simulation studies confirm the proposed estimators' strong performance in finite samples, even under severe contamination. An empirical example is also provided to demonstrate practical relevance. This approach helps to bridge the gap between robustness and computational efficiency, offering a viable alternative to more complex robust methods in the literature. It is particularly useful for large datasets where computational cost is a key concern.

C0898: Testing the null hypothesis of panel cointegration with common factors

Presenter: **Josep Lluís Carrion-i-Silvestre**, Universitat de Barcelona, Spain

Co-authors: Anindya Banerjee

The aim is to address testing the null hypothesis of panel cointegration with cross-section dependence driven by unobserved common factors. The use of the continuously updated estimator proposed in a prior study allows consistent estimation of the cointegrating vector, the common component, and the idiosyncratic component, which establishes the framework to propose a Lagrange multiplier statistic to test the cointegration assumption upon which these results rely.

CO061 Room BCB 212 DYNAMIC MODELS FOR FINANCIAL DATA ANALYSES

Chair: Massimiliano Caporin

C0436: Duration modeling in the presence of zero-duration clusters

Presenter: **Alessandro Morelli**, University of Milan, Italy

Co-authors: Massimiliano Caporin, Eduardo Rossi

The presence of groups of transactions sharing the exact same timestamps in high-frequency financial data poses a significant challenge for duration modeling. This well-known issue is addressed by proposing a novel modeling approach that leverages the predictive information contained in such groups, which is referred to as zero-duration clusters. Using millisecond-level tick data for twelve large-cap, highly liquid U.S. stocks over the last two quarters of 2024, three novel empirical regularities are documented: Durations tend to decrease before a cluster and increase afterward; larger clusters are associated with shorter subsequent durations; clusters exhibit persistence over time, especially when large. Motivated by these findings, an extension of the autoregressive conditional duration (ACD) model is proposed that jointly captures these dynamics through a limited number of additional parameters. An extensive out-of-sample evaluation against standard alternatives, including adaptations of exponential and gamma generalized autoregressive score (GAS) models within the ACD framework, demonstrates that the proposed specification consistently emerges

as the only model within the model confidence set (MCS). It is concluded that while aggregating zero durations enables modeling through point processes, the informational content carried by clusters should nonetheless be exploited through an appropriate modeling approach.

C0493: **Realized co-illiquidity**

Presenter: **Eduardo Rossi**, University of Pavia, Italy

Co-authors: Paolo Santucci de Magistris, Orimar Sauri, Angelo Ranaldo

The purpose is to examine market liquidity, with a focus on co-liquidity and illiquidity, which are crucial for both market participants and policymakers. It addresses the importance of extreme illiquidity events and their risks, particularly in high-frequency trading environments. Using high-frequency data on returns and volume, it builds on a simple structural model and a non-parametric approach inspired by a recent study, which links market volatility, volume, and liquidity. Liquidity measurement is refined by capturing price changes and trading volume dynamics, with a focus on illiquidity due to market frictions. It extends the framework to examine the temporal dynamics of co-(il)liquidity across various assets, providing new metrics to quantify liquidity contagion and systemic risks. This methodology bridges theory with empirical applications, offering practical tools for assessing and managing liquidity risks in diverse financial markets.

C0184: **Investor sentiment and tail risk spillovers**

Presenter: **Petre Caraiani**, Bucharest University of Economic Studies, Romania

Co-authors: Anghel Dan Gabriel

High-frequency data is used to construct measures of sentiment spillovers for financial stocks listed in the United States. A statistically and economically significant response of tail risk spillovers to investor sentiment spillovers is then tested and found. Significant effects are found for both positive and negative sentiment. For the former, this leads to dampening tail risk spillovers, while for the latter, tail risk spillovers are amplified.

C0391: **Forecast reconciliation and multivariate GARCH**

Presenter: **Massimiliano Caporin**, University of Padova, Italy

Co-authors: Emanuele Lopetuso, Daniele Girolimetto

When considering portfolio risk forecasts, for instance, within a risk management perspective, multivariate GARCH models represent a commonly adopted approach. However, if portfolio weights are known, a univariate GARCH model on historically reconstructed portfolio returns can be considered. This represents a framework where forecast reconciliation can be considered as a tool for further improving portfolio risk prediction. By resorting to simulations, the advantage of combining univariate and multivariate portfolio risk forecasts is assessed with the aid of forecast reconciliation techniques. Results suggest relevant advantages when multivariate models are misspecified and the underlying variance is known. However, when a noisy proxy is used, all models, even the misspecified ones, are very close to each other. The analyses are complemented with an empirical example also exploiting the informative content of high-frequency data.

CO043 Room BCB 213 RECENT ADVANCES IN ASYMPTOTIC STATISTICS FOR STOCHASTIC PROCESSES

Chair: Masayuki Uchida

C0512: **Asymptotic inference for skewed stable Ornstein-Uhlenbeck process**

Presenter: **Hiroki Masuda**, University of Tokyo, Japan

Co-authors: Eitaro Kawamo

Asymptotic inference is considered for a non-Gaussian stable Levy process with possibly asymmetric jumps. Based on the infill sampling scheme in a fixed period, the local likelihood asymptotics through a non-diagonal normalizing matrix are presented. Moreover, to estimate the parameters in the driving skewed stable noise, an easy-to-compute empirical moment fitting is proposed. This can serve as an initial estimator for the numerical search of the maximum-likelihood estimator, allowing the bypass of the heavy computational load in optimizing the original log-likelihood function. Simulation results are given to show that the proposed estimator achieves accuracy comparable to that of the maximum likelihood estimator in a much shorter time.

C0819: **Fluid limit of piecewise deterministic Monte Carlo methods**

Presenter: **Kengo Kamatani**, ISM, Japan

Co-authors: Joris Bierkens, Gareth Roberts, Sanket Agrawal

Piecewise deterministic Markov processes (PDMPs) provide the basis for several continuous-time Monte Carlo algorithms, including the bouncy particle sampler (BPS) and the ZigZag process. Their transient phase, meaning the movement from a low-density start to a high-density region, is studied under a convex potential. By combining fluid limits with an averaging decomposition of the generator into fast (nonergodic) and slow parts, the stochastic dynamics are approximated with deterministic ordinary differential equations. For Gaussian targets, BPS and ZigZag require the same early-stage cost as the random-walk Metropolis (RWM). For heavier-tailed targets, well-implemented PDMPs can outperform RWM, and certain event-chain variants achieve dimension-free mixing during the transient phase. These results clarify early behavior and guide the application of PDMP-based Monte Carlo methods to large-scale inference.

C0830: **Drift estimation for rough processes under small noise asymptotic: QMLE approach**

Presenter: **Arnaud Gloter**, Université d'Evry Val d'Essonne, France

Co-authors: Nakahiro Yoshida

The purpose is to consider a process X solution of a stochastic Volterra equation with an unknown parameter θ in the drift function. The Volterra kernel is singular and given by $K(u) = cu^{\alpha-1}$ with $\alpha \in (1/2, 1)$, and it is assumed that the diffusion coefficient is proportional to $\varepsilon \rightarrow 0$. Based on the observation of a discrete sampling with mesh h of the Volterra process, a quasi maximum likelihood estimator is built. The main step is to assess the error arising in the reconstruction of the path of a semi-martingale from the inversion of the Volterra kernel. It is shown that this error decreases as $h^{1/2}$, whatever is the value of α . Then, an explicit contrast function can be introduced, which yields an efficient estimator when $\varepsilon \rightarrow 0$.

C0913: **Deep learning of point processes for modeling high-frequency data**

Presenter: **Nakahiro Yoshida**, University of Tokyo, Japan

Co-authors: Yoshihiro Gytoku, Ioane Muni Toke

Applications of deep neural networks are investigated to a point process having an intensity with mixing covariates processes as input. The generic model includes Cox-type models and marked point processes as well as multivariate point processes. An oracle inequality and a rate of convergence are derived for the prediction error.

CO235 Room BCB M201 NEW CONTRIBUTIONS TO EXTREME VALUE THEORY

Chair: Armelle Guillou

C0417: **Conditional extreme value estimation for dependent time series**

Presenter: **Theodor Henningsen**, University of Copenhagen, Denmark

Co-authors: Martin Bladt, Laurits Glargaard

The consistency and weak convergence of the conditional tail function and conditional Hill estimators is studied under broad dependence assumptions for a heavy-tailed response sequence and a covariate sequence. Consistency is established under alpha-mixing, while asymptotic normality follows from beta-mixing and second-order conditions. A key aspect of the approach is its versatile functional formulation in terms of the condi-

tional tail process. Simulations demonstrate its performance across dependence scenarios. The method is applied to extreme event modeling in the oil industry, revealing distinct tail behaviors under varying conditioning values.

C0428: Corrected inference about the extreme expected shortfall in the general max-domain of attraction (part I)

Presenter: Gilles Stupfler, University of Angers, France

Co-authors: Abdelaati Daouia, Antoine Usseglio-Carleve

The use of the expected shortfall as a solution for various deficiencies of quantiles has gained substantial traction in the field of risk assessment over the last 20 years. Existing approaches to its inference at extreme levels remain limited to distributions that are both heavy-tailed and have a finite second tail moment. This constitutes a strong restriction in areas like finance and environmental science, where the random variable of interest may have a much heavier tail or, at the opposite, may be light-tailed or short-tailed. Under a wider semiparametric extreme value framework, comprehensive asymptotic theory is developed for expected shortfall estimation above extreme quantiles in the class of distributions with finite first tail moment, regardless of whether the underlying extreme value index is positive, negative, or zero. The obtained asymptotic theory is contrasted with what is currently known in the literature.

C0429: Corrected inference about the extreme expected shortfall in the general max-domain of attraction (part II)

Presenter: Antoine Usseglio-Carleve, Avignon Universita, France

Co-authors: Abdelaati Daouia, Gilles Stupfler

The finite-sample applicability of asymptotic theory for expected shortfall estimation above extreme quantiles in the class of distributions with finite first tail moment, regardless of whether the underlying extreme value index is positive, negative, or zero, is discussed and found to often fail to yield confidence intervals whose empirical coverage probability matches nominal coverage. By relying on the moment estimators of the scale and shape extreme value parameters, as well as on a fine understanding of the dependence structure between these estimators and intermediate expected shortfall estimators, corrected asymptotic confidence intervals are constructed whose finite-sample coverage is found to be close to the nominal level on simulated data. The usefulness of the construction is illustrated on two sets of financial loss returns and flood insurance claims data.

C0551: Asymptotic theory for the likelihood-based block maxima method in time series

Presenter: Simone Padoan, Bocconi University, Italy

Co-authors: David Carl, Stefano Rizzelli

An asymptotic framework is developed for likelihood-based inference in the block maxima (BM) method for stationary time series. While Bayesian inference under the BM approach has been widely studied in the independence setting, no asymptotic theory currently exists for time series. Further results are needed to establish that the BM method can be applied with the kind of dependent time series models relevant to applied fields. To address this gap, a comprehensive likelihood theory is first established for the misspecified generalized extreme value (GEV) model under serial dependence. Building on this foundation, the asymptotic theory of Bayesian inference is developed for the GEV parameters, the extremal index, T -time-horizon return levels, and extreme quantiles (value at risk). For inference on the extremal index, an adjusted posterior distribution is proposed that corrects for poor coverage exhibited by a naive Bayesian approach. Simulations show excellent inferential performances for the proposed methodology.

CO086 Room BCB 308 MACHINE LEARNING AND HIGH-DIMENSIONAL DATA

Chair: Eoghan O'Neill

C0801: Model-free identification in ill-posed regression

Presenter: Gianluca Finocchio, University of Vienna, Austria

Co-authors: Tatyana Krivobokova

The problem of parsimonious parameter identification in possibly high-dimensional linear regression with highly correlated features is addressed. This problem is formalized as the estimation of the best, in a certain sense, linear combinations of the features that are relevant to the response variable. Importantly, the dependence between the features and the response is allowed to be arbitrary. Necessary and sufficient conditions for such parsimonious identification – referred to as statistical interpretability – are established for a broad class of linear dimensionality reduction algorithms. Sharp bounds on their estimation errors, with high probability, are derived. To best of knowledge, this is the first formal framework that enables the definition and assessment of the interpretability of a broad class of algorithms. The results are specifically applied to methods based on sparse regression, unsupervised projection, and sufficient reduction. The implications of employing such methods for prediction problems are discussed in the context of the prolific literature on overparametrized methods in the regime of benign overfitting.

C0827: Nonlinear autoregressive models for functional time series with Bayesian additive regression trees

Presenter: Eoghan O'Neill, Erasmus University Rotterdam, Netherlands

Co-authors: Maria Grith, Anastasija Tetereva, Jiahao Cao, Guanyu Hu

The purpose is to introduce Bayesian additive regression tree models for function-on-function regression. The outcome function is modelled as a linear combination of data-adaptive basis functions. The coefficients of basis functions are determined by sums of trees that can split on both scalar and functional variables, including the lag of the dependent variable. Splitting rules for functions are defined by inner products between functional covariates and linear combinations of fixed basis functions, distinct from the aforementioned data-adaptive basis. A prior is considered on functional splits, allowing Markov chain Monte Carlo tree samples to adapt to the data by placing higher probability on selecting a subset of relevant basis functions for a splitting rule. The forecasting performance of the method is evaluated in an application to option pricing implied volatility surface data.

C1061: Common factors in currency characteristics

Presenter: Dennis Umlandt, University of Innsbruck, Austria

Co-authors: Moritz Dauber

The factor structure of currency characteristics are studied by employing a three-dimensional tensor factor model that simultaneously captures the variation in characteristics of the G10 currencies over time. It is shown that factor-mimicking portfolios derived from these common factors in currency characteristics are able to price individual currency returns better than standard factor models derived from univariate sorts on the same characteristics. The variation in currency characteristics can be well captured by a parsimonious two-factor model, where the first factor closely resembles the carry trade and the second factor acts as a hedge against carry crash risk, which is composed of signals from FX momentum, FX value, and the term spread. A potential third factor, which dynamically weights several characteristics, incrementally improves the fit of the total variation, but has a high Sharpe ratio.

C1068: Learning continuous-time network dynamics via network and sheaf SDEs

Presenter: Francesco Iafate, University of Hamburg, Germany

Learning and inference is investigated for continuous-time network dynamics and their topological generalizations. First, a class of network stochastic differential equations (N-SDEs) is proposed, in which every node of a directed graph follows an SDE driven by: (i) a momentum term capturing the nodes own past, (ii) a network effect that aggregates feedback from its neighbors, and (iii) stochastic volatility from Brownian noise. Non-asymptotic error bounds are provided for joint parameter estimation and graph recovery in a high-frequency (ergodic) observation scheme, and design an adaptive-lasso routine that learns the graph when it is unknown. The possibility of exploiting graph sparsity and modeling oriented

edges makes the model attractive for causal inference. Motivated by the growing success of topological methods in data science, N-SDEs are then extended to cellular sheaves, obtaining the Sheaf-SDEs. A cellular sheaf equips every node and edge with a vector space and restriction maps; embedding SDE dynamics in this structure lets multivariate signals evolve under local self-interaction, cross-variable coupling, and higher-order influences that propagate through the sheaf, capturing relationships that ordinary graph models cannot express. Synthetic experiments validate exact structural recovery, while a case study on pollutant diffusion over spatial networks illustrates how the framework allows for powerful location-dependent models without sacrificing interpretability.

CO231 Room BCB 309 TOPOLOGICAL AND GEOMETRIC STATISTICAL ANALYSIS
Chair: Carlos Soto
C0352: Persistence diagrams, use and limitations

Presenter: **Olympio Hacquard**, Kyoto University, Japan

Persistence diagram is a central tool in topological data analysis, aiming at characterizing all topological information present in a dataset. It is shown how to use it to perform a regression task by providing a topological characterization of functional noise. Given its high level of detail, persistence diagrams suffer from certain computational drawbacks, limiting their practical use. Some alternative topological descriptors are presented based on Euler characteristic computation that provide a high classification performance at a much reduced computational cost.

C0415: Topological analysis for detecting anomalies in time series

Presenter: **Clement Levrard**, Universita de Rennes, France

A recent methodology is exposed based on the field of topological data analysis for detecting anomalies in multivariate time series, which aims to detect global changes in the dependency structure between channels. This approach is lean enough to handle large-scale datasets, and extensive numerical experiments back the intuition that it is more suitable for detecting global changes of correlation structures than existing methods. If time allows, some theoretical guarantees will also be presented for this method.

C0495: Wormlike chain model and spherical spectral gap

Presenter: **Ho Yun**, EPFL, Switzerland

Co-authors: Almond Stoecker, Victor Panaretos

How can the stiffness of a DNA molecule or a protein filament be precisely measured? The wormlike chain model, a foundational framework in structural biology, offers a mathematical description of these biopolymers. The central parameter of this model is called the persistence length, quantifying their bending stiffness. A connection is revealed between this quantity and the spectral gap of a diffusion operator on the sphere. Using the theory of spherical harmonics, this relationship naturally leads to a classic model in directional statistics to describe polymer conformations. Building on this observation, generalized statistical models are introduced to describe DNA sequential data and discuss compelling open questions at the intersection of statistics and molecular biology.

C0903: Trajectory inference with varifold distances

Presenter: **Elodie Maignant**, Zuse Institute Berlin, Germany

Co-authors: Christoph von Tycowicz, Tim Conrad

The focus is on a tree inference problem motivated by the problem, known as trajectory inference, in single-cell genomics of reordering a population of cells sampled from a dynamic process according to their progression in the process. If the process is differentiation, this amounts to reconstructing the corresponding differentiation tree. One way of doing this in practice is to estimate the shortest-path distance between nodes based on cell similarities observed in sequencing data. Recent sequencing techniques make it possible to measure two types of data: Gene expression levels and RNA velocity, a vector that predicts changes in gene expression. The data then consist of a discrete vector field on a (subset of a) Euclidean space of dimension equal to the number of genes under consideration. By integrating this velocity field, the evolution of gene expression levels is traced in each single cell from some initial stage to its current stage, and using varifold distances between the curves thus obtained, a similarity measure is defined between nodes which is proven to approximate the shortest-path distance in a tree that is isomorphic to the differentiation tree.

CO211 Room BCB 313 RECENT ADVANCES IN SURVIVAL DATA ANALYSIS
Chair: Chi Hyun Lee
C0275: Variable selection in ultra-high dimensional feature space for the Cox model with interval-censored data

Presenter: **Daewoo Pak**, Yonsei University, Korea, South

The aim is to present a series of variable selection methods for the Cox model with interval-censored data, tailored for ultra-high-dimensional settings where the number of covariates may grow exponentially with the sample size. The methods select covariates via a penalized nonparametric maximum likelihood estimation with some popular penalty functions, including lasso, adaptive lasso, SCAD, and MCP. It is proven that the penalized variable selection methods with folded concave penalties or adaptive lasso penalty enjoy the oracle property. Extensive numerical experiments show that the proposed methods have satisfactory empirical performance under various scenarios. The utility of the methods is illustrated through an application to a genome-wide association study of age to early childhood caries.

C0604: Self-consistent equation-guided neural networks for censored time-to-event data

Presenter: **Sehwan Kim**, Ewha Womans University, Korea, South

Co-authors: Rui Wang, Wenbin Lu

In survival analysis, estimating the conditional survival function given predictors is often of interest. There is a growing trend in the development of deep learning methods for analyzing censored time-to-event data, especially when dealing with high-dimensional predictors that are complexly interrelated. Many existing deep learning approaches for estimating the conditional survival functions extend the Cox regression models by replacing the linear function of predictor effects by a shallow feed-forward neural network while maintaining the proportional hazards assumption. Their implementation can be computationally intensive due to the use of the full dataset at each iteration, because the use of batch data may distort the at-risk set of the partial likelihood function. To overcome these limitations, a novel deep learning approach is proposed for non-parametric estimation of the conditional survival functions using the generative adversarial networks, leveraging self-consistent equations. The proposed method is model-free and does not require any parametric assumptions on the structure of the conditional survival function. The convergence rate of the proposed estimator of the conditional survival function is established. In addition, the performance of the proposed method is evaluated through simulation studies, and its application on a real-world dataset is demonstrated.

C0787: Imputation-based q-learning with regression trees for censored survival data

Presenter: **Youngjoo Cho**, Konkuk University, Korea, South

Co-authors: Xue Yang, Abdus Wahed, Yu Cheng

Dynamic treatment regimes (DTRs) are sets of decision rules that guide individualized, time-varying treatments in multistage therapy. In the presence of right-censored survival data, many methods have been proposed to determine optimal DTRs that maximize the expected overall survival time. Q-learning is a commonly used and straightforward reinforcement learning algorithm for this purpose. However, it is sensitive to model misspecification, a well-known limitation. To address this issue, tree-assisted imputation-based Q-learning (TAI Q-learning) is proposed. In this method, double-robust survival trees and tree ensembles are used to estimate the optimal treatment rules at each stage, while multivariate imputation by chained equations (MICE) is employed to predict the optimal survival time. Through extensive simulation studies, the performance of traditional Cox proportional hazards (Cox PH) and accelerated failure-time (AFT) models with nonparametric tree-based methods in the optimization step are

compared, and hot-deck multiple imputation is compared with MICE in the imputation step. The simulation results show that MICE is easier to implement and often outperforms hot-deck imputation. Moreover, in multilevel treatment scenarios, tree-based methods outperform standard Cox PH and AFT models in estimating optimal treatment rules.

C1104: Rank estimation of monotone individualized treatment regimes for survival outcomes

Presenter: Taehwa Choi, Sungshin Women's University, Korea, South
Co-authors: Seohyeon Park, Hyeonseok Oh, Zhezhen Jin, Sangbum Choi

An optimal individualized treatment regimen (ITR) recommends tailored treatment decisions based on patients' genetics and demographic information, thereby providing greater clinical benefits compared to traditional clinical trials that randomly assign binary treatments to patients without considering their clinical status. While many studies have been conducted over several decades to identify the optimal ITR, the majority of them have covered simpler models that cannot be generalized to more complex models, such as the single-index model. The aim is to propose an inferential procedure for the surrogate maximum rank correlation (SMRC) to determine the optimal ITR under the single-index transformation model. Unlike the existing methods, the proposed approach provides greater flexibility, avoids strict model assumptions, and eliminates the need for complex computations while maintaining statistical accountability. Furthermore, an efficient variance estimation procedure is developed, using induced smoothing. Large sample properties are investigated, and various numerical examples demonstrate the usefulness of the method.

CO274 Room BCB 402 VARIABLE IMPORTANCE AND STATISTICAL LEARNING IN ENVIRONMENT AND ENERGY Chair: Consuelo Nava

C0492: Understanding dependencies in compositional data through graphical models

Presenter: Agnese Maria Di Brisco, University of Piemonte Orientale, Italy
Co-authors: Roberto Ascari, Federica Nicolussi, Anna Maria Fiori

A methodology is introduced for analyzing (in)dependencies in compositional data using graphical models. Compositional data are first mapped to an unconstrained space, where Gaussian graphical models serve as a starting point for representing dependency structures. A novel subclass of these models is defined and proposed, tailored to respect the specific constraints of compositional data through block-diagonal covariance structures. An extension to the non-Gaussian setting is introduced, while preserving the compositional structure of the data. Estimation is based on block-diagonal covariance matrices, with model selection guided by a penalized likelihood criterion and cross-validation. The methodology is validated through applications to both simulated and real-world data, demonstrating its ability to reveal meaningful relationships in compositional contexts.

C0553: An exact game-theoretic variable importance index for generalized additive models

Presenter: Amir Khorrami Chokami, University of Cagliari, Italy
Co-authors: Giovanni Rabitti

Generalized additive models (GAMs) are a widely adopted tool in statistical modeling. The problem of assessing variable importance in GAMs is addressed by introducing a variance allocation approach based on the Shapley value. A closed-form expression is derived for this importance index, allowing for efficient computation in high-dimensional settings and under general dependence structures. The practical implication is discussed that when a variable's importance is negligible, it can be safely eliminated from the GAM, simplifying the model. The case studies show that the Shapley values offer more informative insights than p-values in terms of ranking the importance of variables.

C0885: Greener strategies, stronger results? Environmental sustainability and firm performance in Italy

Presenter: Luca Patelli, University of Bergamo, Italy
Co-authors: Peter Cincinelli, Daniele Toninelli, Giovanna Zanotti

In the current context of environmental challenges, international institutions have established targets with the aim of mitigating the impact of climate change. The aim, focused on Italian firms, is to investigate whether sustainability, and more specifically the environmental dimension, affects the probability of observing firms in a state of financial distress. The probability of distress is estimated by considering financial and accounting variables, as well as variables derived from non-financial disclosures and good sustainability practices; variables available in external ESG databases are also taken into account. To predict the probability of distress, a statistical machine learning framework is relied on. These methods enable flexible handling of a large number of predictors of different types. In order to overtake the black-box nature of several machine learning algorithms, a set of explainable AI solutions is implemented. This approach allows assessing the variables' importance as well as to enhance the model's explainability, thereby facilitating a deeper and clearer comprehension of the prediction mechanisms and of the obtained results.

C0960: Bootstrapping time series of renewable energy sources for environmental and energy studies

Presenter: Giulia Marcon, Università degli studi di Palermo, Italy

The increasing integration of renewable energy sources (RES) into power systems requires effective methods to handle their intrinsic variability, especially for small-scale systems like microgrids. These systems, often operating in islanded mode due to main grid disruptions, demand robust scenario-based analyses to ensure operational stability and energy resilience. A methodology is presented for generating synthetic time series of key meteorological variables (wind speed, solar irradiance, and temperature) used to model RES variability under diverse environmental conditions. The proposed approach employs a time series bootstrap technique, which preserves both short-term autocorrelations and seasonal patterns, while avoiding strong parametric assumptions. The resulting synthetic datasets enable the creation of multiple stochastic realizations, supporting environmental and energy studies focused on system reliability, performance evaluation, and risk assessment. These scenarios are particularly valuable for microgrid planning and operation, especially in isolated or fault-prone contexts. The method is being integrated into a broader framework for optimal energy storage utilization and microgrid reliability analysis under uncertainty.

CO161 Room BCB 405 BAYESIAN STRUCTURE LEARNING IN GRAPHICAL MODELS

Chair: Reza Mohammadi

C0323: Exact Bayesian computation for large Gaussian graphical models

Presenter: David Rossell, Universitat Pompeu Fabra, Spain
Co-authors: Jack Jewson, Deborah Sulem

Bayesian methods possess appealing properties for quantifying the uncertainty associated with learning the dependence structure in a graphical model from data, as well as the uncertainty in parameter estimates. Computational bottlenecks limited their application when the number of variables is large, which prompted the use of pseudo-Bayesian approaches. Computational algorithms are proposed for exact Bayesian inference that provably scale well to high-dimensional settings, when the data-generating precision matrix is sparse. The framework is based on discrete spike-and-slab priors under which, by exploiting sparsity, spectral gaps and the per-iteration cost can be bound. MCMC algorithms are proposed that allow row-wise updates of the precision matrix, either using standard local proposals (e.g., Gibbs, birth-death-swap, LIT) or a novel global proposal that may add/remove multiple edges in one iteration. Examples show that the methods extend the applicability of exact Bayesian inference from roughly 100 to roughly 1,000 variables (equivalently, from 5,000 edges to 500,000 edges, or from 2^5000 models to 2^500000 models).

C0358: Fast MCMC chains for Bayesian posterior inference in graphical models

Presenter: Lucas Vogels, University of Amsterdam, Netherlands
Co-authors: Reza Mohammadi, Ilker Birbil, Marit Schoonhoven

The purpose is to treat a problem as old as statistics: Discovering the unknown parameters of a distribution. A Bayesian framework is used for

this. That means the posterior distribution is set out to be discovered. Markov chain Monte Carlo (MCMC) methods are a popular tool for this. Traditional MCMC strategies, such as reversible jump or birth-death algorithms, are still popular, despite suffering from a slow exploration of the parameter space. An alternative is offered by approximating the detailed balance condition, creating an algorithm that can traverse the entire parameter space in a single iteration. The power of the algorithm is shown in the field of Gaussian graphical models, Ising models, and feature selection, but it is noted that its applicability reaches beyond those examples. In fact, it is applicable to most Bayesian posterior inference.

C0461: Efficient sampling for Bayesian networks

Presenter: **Jack Kuipers**, ETH Zurich, Switzerland

Bayesian networks are probabilistic graphical models widely employed to understand dependencies in high-dimensional data, and even to facilitate causal discovery. Learning the underlying network structure, which is encoded as a directed acyclic graph (DAG), is highly challenging, mainly due to the vast number of possible networks in combination with the acyclicity constraint, and a wide plethora of algorithms have been developed for this task. Efforts have focused on two fronts: Constraint-based methods that perform conditional independence tests to exclude edges and score, and search approaches that explore the DAG space with greedy or MCMC schemes. These two fields are synthesized in a novel hybrid method, which reduces the complexity of Bayesian MCMC approaches to that of a constraint-based method. This enables full Bayesian model averaging for much larger Bayesian networks and offers significant improvements in structure learning. The application of Bayesian network learning to patient stratification in myeloid malignancies is further discussed.

C0647: A spatial autoregressive graphical model

Presenter: **Pariya Behrouzi**, Wageningen University and Research, Netherlands

Co-authors: Sjoerd Hermes, Joost van Heerwaarden

There is a notable gap in the statistical literature for methods capable of modeling asymmetric multivariate spatial effects, particularly in settings where spatial relationships vary across categorical labels. In such scenarios, observations at a location arise from both within- and between-location effects, with the latter often exhibiting asymmetry due to heterogeneous interactions between different location types. A novel Bayesian spatial graphical model is proposed that integrates multivariate spatial autoregressive structures with Gaussian graphical models. This integration allows capturing asymmetric spatial dependencies that are modulated by a categorical feature at each location. These feature-dependent spatial effects relax the usual symmetry assumptions commonly imposed in spatial models. However, the added flexibility comes with a trade-off: Spatial effects are not identifiable without prior knowledge of the system or additional parameter constraints. The model's performance is evaluated via simulation studies, and its practical utility is demonstrated on an intercropping dataset, where asymmetric spatial effects naturally arise from interactions between different crop types. The proposed method is implemented in the R package SAGM.

CO273 Room BCB 406 BAYESIAN ANALYSIS OF NETWORK DATA

Chair: Sirio Legramanti

C0809: Bayesian community detection in assortative stochastic block model with an unknown number of communities

Presenter: **Martina Amongero**, University of Torino, Italy

Co-authors: Pierpaolo De Biasi

Available data in the form of networks is gaining increasing attention in modern research; examples include social networks, biological networks, and many others. A statistical problem of key interest is community detection, that is, to divide the nodes into strongly connected clusters with relatively weak cross-cluster connections. The stochastic block model (SBM) provides a well-suited generative process for explaining the formation of communities. By leveraging the SBM, researchers gain insights into the underlying structure of a network, uncovering interaction patterns that may not be apparent from raw data, and exploring network properties such as assortativity. In particular, an SBM is assortative when the probability of a connection is higher when nodes belong to the same rather than to different clusters. A recent line of work uses Bayesian non-parametric methods for the recovery of communities in classical SBM by placing a prior distribution on the number of clusters and estimating cluster assignments with collapsed Gibbs samplers. However, the development of an efficient Gibbs sampler for assortative SBM is still an open problem. The aim is to enforce the assortative property through the prior and to study its effect on community detection by implementing a Gibbs sampler for posterior inference. A simulation study based on a class of benchmark datasets is conducted to demonstrate the benefits of the assortative case over the standard one.

C0457: Dynamic network models with time-varying nodes

Presenter: **Luca Gherardini**, University of Klagenfurt, Austria

Co-authors: Monia Lupporelli, Mauro Bernardi

Dynamic networks offer a versatile framework for examining temporal dependencies among statistical units, yielding insights into the stochastic mechanisms that drive interactions over time. An often-overlooked complexity - population dynamics - is addressed, whereby individuals may enter or exit the network during the observation window, potentially biasing inferential outcomes if unaccounted for. Focusing on binary temporal networks with evolving topologies, it distinguishes between two scenarios: (i) fully observed network evolution and (ii) partially observed topology requiring inferential reconstruction. In both contexts, a hierarchical mixed-effects model that jointly characterizes uncertainty in node-set evolution and dyadic links is proposed. To enable scalable inference, a Bayesian conjugate algorithm grounded in Polya-Gamma data augmentation is developed. The principal inferential properties of the model are derived, establishing its theoretical validity. Through simulation experiments and an application to social interaction data from a school setting, the approach is demonstrated to not only mitigate biases associated with time-varying populations but also enhance the interpretability of dynamic network patterns.

C0745: A zero-inflated Poisson latent position cluster model

Presenter: **Riccardo Rastelli**, University College Dublin, Ireland

Co-authors: Chaoyi Lu, Nial Friel

The latent position model is a popular approach for the statistical analysis of network data. A key aspect of this model is that it assigns nodes to random positions in a latent space, and the probability of an interaction between each pair of nodes is determined by their distance, allowing researchers to visualize nuanced structures via a latent embedding of the graph. Missing data is a common issue in statistical network analysis, often leading to an excess of observed zeros in interaction data. The focus is on non-negatively weighted social networks. By treating missing data as 'unusual' zero interactions, a combination of the zero-inflated Poisson distribution is proposed with the latent position cluster model. The framework extends the latent position model to accommodate the clustering of individuals and to simultaneously model weighted interactions, handle missing data, perform clustering, and produce three-dimensional visualizations of real networks. Statistical inference is based on a partially collapsed Markov chain Monte Carlo algorithm, which involves a new truncated absorb-eject move, and selects the number of groups automatically by leveraging a mixture of finite mixtures framework.

C0510: Dependent stochastic block models for sequences of directed networks with application to causes of death co-occurrences

Presenter: **Giovanni Romano**, Bocconi University, Italy

Co-authors: Cristian Castiglione, Daniele Durante

A new Bayesian model is developed for inferring changes in the stochastic block structures within a sequence of weighted and directed networks indexed by a predictor. This model, named dynamic stochastic block model (dSBM), is originally motivated by the demographic problem of learning age-specific group structures in directed networks encoding co-occurrences among underlying and multiple death causes for different age groups in a population. To this goal, state-of-the-art stochastic block models are substantially generalized to account for (i) edges with categorical

weights, (ii) two separate node partitions for their different roles in a directed network, i.e., sender/receiver, and (iii) sequences of networks indexed by an ordered predictor, such as age or time. Results in the causes-of-death networks analysis unveil interesting and yet-unexplored patterns in the composition, evolution, and interactions of causes-of-death clusters, with relevant analytic and policy implications. The novel formulation we propose provides a powerful model to infer non-trivial changes in grouping structures governing sequences of dynamic directed networks, beyond the demographic application.

CO115 Room BCB 408 BAYESIAN SEMI- AND NON-PARAMETRIC METHODS III
Chair: Beatrice Franzolini
C0634: A Bayesian nonparametric framework for dynamic item-response theory

Presenter: **Maria De Iorio**, National University of Singapore, Singapore

Item-response theory (IRT) is widely used for the statistical analysis of questionnaire data, allowing for the differentiation of respondent profiles and the characterization of questionnaire items through interpretable parameters. A Bayesian semiparametric extension of IRT is proposed that introduces temporal dependence across repeated questionnaire administrations, accommodates repeated measurements, and jointly models responses from related subject groups (e.g., mothers and children) to enable information sharing across hierarchies. The framework further incorporates covariate information, allows for the joint modeling of questionnaire data with other longitudinal markers, and supports clustering of subjects based on their latent response profiles. The approach is built on Bayesian nonparametric priors: The Dirichlet process and the normalized generalized gamma process, facilitating the identification of clinically meaningful subgroups. The utility of the proposed methodology is demonstrated through the analysis of longitudinal psychometric questionnaire data from the Singaporean GUSTO cohort study. These data are collected from mothers and their children, aiming to investigate how various factors influence growth trajectories, developmental outcomes, and mental health.

C1170: Spectral decomposition-assisted multi-study factor analysis

Presenter: **Niccolo Anceschi**, Duke University, United States

Co-authors: Lorenzo Mauri, David Dunson

The focus is on covariance estimation for multi-study data. Popular approaches employ factor-analytic terms with shared and study-specific loadings that decompose the variance into (i) a shared low-rank component, (ii) study-specific low-rank components, and (iii) a diagonal term capturing idiosyncratic variability. The proposed methodology estimates the latent factors via spectral decompositions and infers the factor loadings via surrogate regression tasks, avoiding identifiability and computational issues of existing alternatives. Reliably inferring shared vs study-specific components requires novel developments that are of independent interest. The approximation error decreases as the sample size and the data dimension diverge, formalizing a blessing of dimensionality. Conditionally on the factors, loadings and residual error variances are inferred via conjugate normal-inverse gamma priors. The conditional posterior distribution of factor loadings has a simple product form across outcomes, facilitating parallelization. Favorable asymptotic properties are shown, including central limit theorems for point estimators and posterior contraction, and excellent empirical performance in simulations. The methods are applied to integrate three studies on gene associations among immune cells.

C0841: Autocompound random measures

Presenter: **Riccardo Corradin**, University of Milano-Biocca, Italy

Co-authors: Fabrizio Leisen

The aim is to introduce a class of time-dependent nonparametric models, suited for scenarios involving populations observed at distinct discrete times. Starting with an ancestral random measure, it is possible to define a sequence of time-specific random measures by acting on their intensity functions, in the spirit of compound random measures. The resulting family of models exhibits desirable properties, including mathematical tractability, simple expressions for its main summaries, and a closed-form representation of the joint posterior distribution at distinct observed times. Such a model can then be normalized and used as a building block for dynamic population studies, defining tractable species sampling models that evolve over time, or convoluted with a kernel function to obtain time-dependent mixture models.

C0689: Oh SnapMMD! Forecasting stochastic dynamics beyond the Schrodinger bridge's end

Presenter: **Renato Berlinghieri**, Massachusetts Institute of Technology, United States

Scientists often need to predict system behavior beyond the time window covered by snapshot data governed by latent stochastic dynamics. In single-cell mRNA profiling, for instance, transcriptional states are observed from different replicates at discrete times, but each measurement destroys the cell, so an individual cell's full trajectory is never seen. Yet, researchers want to forecast outcomes (e.g., stem-cell differentiation) from early measurements. Schrodinger-bridge (SB) methods can interpolate between snapshots, but existing approaches either follow a predefined reference dynamic chosen before observing data or assume a fixed, state-independent volatility by minimizing a Kullback-Leibler divergence, both of which can hurt forecasting accuracy. SnapMMD is introduced, a framework that learns latent dynamics by jointly fitting the distribution of observed states and their sampling times using a maximum mean discrepancy (MMD) loss. Unlike prior work, SnapMMD infers unknown, state-dependent volatility directly from the data. In experiments on synthetic and real datasets, SnapMMD delivers more accurate forecasts than SB baselines. It naturally handles missing or partial state observations and provides an interpretable R^2 -style diagnostic of fit quality. Furthermore, SnapMMD matches or exceeds state-of-the-art methods in interpolation tasks and velocity-field reconstruction across all tested scenarios.

CO074 Room BCB 409 SPATIO-TEMPORAL MODELING AND INFERENCE
Chair: Fabian Mies
C0525: The cost of ignoring space: Bias and overconfidence in model calibration

Presenter: **Michele Nguyen**, Nanyang Technological University, Singapore

Co-authors: Maricar Rabonza, David Lallemand

Data with spatial information are collected through active and passive methods such as field surveys and sensor networks. These are used for calibrating models in fields like environmental science, ecology, urban planning, and public health. Despite Tobler's first law of geography, which states that nearby observations tend to be more similar, spatial dependence is sometimes ignored during modeling or omitted by using cost functions that treat observations as independent. An issue in model calibration with spatial data is highlighted: When spatial dependence is present, one may be overconfident in estimated parameters, and in worst cases, face systematic bias. This is evident when there is spatial clustering of observations, and extreme cases see complications from unbalanced data. This is demonstrated through simulation experiments with a simple quadratic model where one dataset has equidistant observations while the other has clustered observations and an isolated point. Next, spatial weights are developed based on spatial conditional information, and it is shown that weighting serves as a middle ground between explicit spatial modeling and omitting spatial dependence. This is illustrated using case studies of heavy metal soil contamination and malaria prevalence. The need for further development of computationally efficient methods to address the complex interplay between clustering, spatial dependence, and model parameters is highlighted.

C0640: Inference in nonstationary random fields and rates of convergence of threshold variance estimation under sparsity

Presenter: **Ansgar Steland**, RWTH Aachen University, Germany

Inference in nonlinear random fields is studied under nonstationarity. The results include a functional central limit theorem for the smoothed spatial partial sum process indexed by classes of sets satisfying a weak entropy condition satisfied by various relevant classes arising in statistics and machine learning. For statistical applications, consistent estimation of the asymptotic variance is needed. The aim is to study (soft-) threshold estimation under a wide class of fields under sparse dependence by establishing nonasymptotic rates of convergence. The results reveal that thresholding is superior under sublinear growth of the non-vanishing spatial covariances over increasing rectangles. Those results also cover

estimation from residuals. Applications to hypothesis testing in images, e.g., to detect cancer, regression models with external regressors, and sparse convolutional network layers are discussed.

C1162: A random graph-based autoregressive model for networked time series

Presenter: **Weichi Wu**, Tsinghua University, China

Co-authors: Chenlei Leng

Contemporary time series data often feature objects connected by a social network, which naturally induces temporal dependence among connected neighbors. The network vector autoregressive model is useful for describing the influence of linked neighbors, while its recent generalizations aim to separate influence and homophily. Existing approaches, however, require either correct specification of a time series model, accurate estimation of a network model, or both, and rely exclusively on least squares for parameter estimation. A new autoregressive model, incorporating a flexible form for latent variables used to depict homophily, is proposed. A first-order differencing method is developed for the estimation of influence, requiring only the influence part of the model to be correctly specified. When the homophily part is correctly specified, admitting a semiparametric form, we leverage and generalize the recent notion of neighbor smoothing for parameter estimation, bypassing the need to specify the generative mechanism of the network. A new theory is developed to show that all the estimated parameters are consistent and asymptotically normal. The efficacy of the approach is confirmed via extensive simulations and an analysis of a social media dataset.

C1399: A universal method for statistical inference of low-frequency time series

Presenter: **Alexis Derumigny**, Delft University of Technology, Netherlands

Co-authors: Fernando De Diego Avila, Fang Fang

Statistical inference for low-frequency time series is challenging due to limited non-overlapping observations, as commonly encountered in fields like financial risk management and weather forecasting. Traditional direct methods relying solely on low-frequency data often produce inaccurate estimates, while estimations based on overlapping observations have biases caused by autocorrelation. A novel simulation-based method is proposed for inferring low-frequency time series. It involves three steps: estimating the distribution of the corresponding higher-frequency process, simulating paths from this distribution, and generating a large dataset of aggregated low-frequency observations to enable accurate and robust statistical inference. Theoretical results are also given on the asymptotic properties of the estimators. The superiority of the proposed method over direct methods is verified via a comprehensive simulation study.

CC371 Room BCB 209 PORTFOLIO OPTIMIZATION

Chair: Yasuhiro Omori

C1003: Income-based optimal portfolio choice: A new analysis

Presenter: **Seyoung Park**, University of Nottingham, United Kingdom

Co-authors: Alain Bensoussan

A new income-based optimal portfolio choice framework is developed. The new state variable and the new risk control representing, respectively, total wealth (financial wealth plus human capital) and the total risk exposure of total wealth to both the stock and labor markets are introduced. In particular, total wealth plays an important direct input variable in the agent's optimal risk control. A certain threshold of wealth is found below which the agent's risk-taking decreases as unspanned income risk increases, and above which the risk-taking rises as the unspanned risk increases. A general-equilibrium analysis demonstrates that the market volatility pattern inherits the agent's different risk-taking behaviors over levels of wealth, so the increased volatility occurs as a rational response to the agent's aggressive bidding for the total risk exposure to both the stock and labor markets when wealth is large.

C1268: Risk budgeting portfolios for general risk measures

Presenter: **Pierpaolo Uberti**, University of Milano-Bicocca, Italy

Co-authors: Claudia Fassino

Given a reference risk measure, the risk budgeting portfolio is defined as the allocation in which each asset contributes a predetermined share to the total risk. This research introduces a novel approach, distinct from those previously proposed in the literature, for the computation of the risk budgeting portfolio. Assuming the existence and uniqueness of the risk budgeting portfolio, a Cauchy sequence is constructed within the simplex, whose limit coincides with the desired portfolio. This construction enables the development of an efficient algorithm that circumvents the need to solve auxiliary optimization problems, which are often computationally demanding and less transparent in the context of decision theory. The proposed algorithm is compared with standard optimization-based methods commonly used in the literature. From a theoretical standpoint, the Cauchy sequence is shown to induce a mapping for which the risk budgeting portfolio represents a fixed point. Consequently, sufficient conditions for the existence and uniqueness of this fixed point can be applied. The methodology is formulated for a general class of risk measures and is illustrated in detail for the case of standard deviation.

C1281: A regularized regression approach to global minimum variance allocation

Presenter: **Timm Pfeil**, University of Innsbruck, Austria

Co-authors: Dennis Umlandt

A regularized regression approach for constructing the global minimum variance portfolio (GMVP) is proposed in large equity universes. Rather than shrinking the return covariance matrix, allocation weights formulated as regression coefficients are directly penalized. This regression-based formulation reduces the number of parameters to be estimated and ensures well-posed solutions even when the number of assets exceeds the sample length ($N > T$). The method is applied to estimate GMVP weights for S&P 500 constituents from 1957 to 2022. Relative to static and dynamic benchmark approaches, the penalized allocations consistently reduce realized out-of-sample variance and improve Sharpe ratios while maintaining comparable or higher average returns. Lasso and elastic net produce sparse, low-turnover allocations, whereas ridge yields stable, diversified exposures. The results are robust across window lengths, index constructions, and crisis subsamples, while preserving the GMVP interpretation, in which regression coefficients map directly to portfolio weights.

CC421 Room BCB M202 TAIL-RISK MODELING

Chair: Jose Olmo

C0791: Sentiment regimes for predicting tail-risk: Tree-structured conditional autoregressive value-at-risk

Presenter: **Christoph Hirt**, University of St Gallen, Switzerland

The empirical importance of sentiment and attention measures is evaluated for conditional quantile forecasting of U.S. equity return processes. To this end, Tree-CAViaR is proposed, a data-driven approach to select among threshold nonlinear CAViaR specifications. Similar to tree-structured models within the GARCH and HAR families, Tree-CAViaR offers a fully data-driven approach to identify structural changes in the quantile process by selecting threshold-regime CAViaR models through binary splitting. Results show that Tree-CAViaR effectively detects distinct tail-risk regimes, defined by exogenous predictors such as the CBOE volatility index (VIX), leading to improved out-of-sample forecasting performance.

C1083: Comparing value at risk and expected shortfall estimation using LSTM and EGARCH family models

Presenter: **Shujie Li**, Paderborn University, Germany

The forecasting performance of value at risk and expected shortfall is examined across several models from the EGARCH family, including the newly introduced modulus log-GARCH, a modified EGARCH, and their long memory extensions, as well as a recurrent neural network based on long short term memory architecture. To assess model performance, eight stock indices from diverse international markets are analyzed. The

models are evaluated using three different back-testing approaches and a model selection criterion, namely the weighted absolute deviation. The results indicate that the selected indices exhibit heavy tails and asymmetry. Therefore, in general, models estimated with skewed distributions perform better than their symmetric counterparts. In most cases, the long short term memory model is selected as the top-performing model. Nevertheless, several models from the EGARCH family remain strong competitors, especially for asymmetric distributions, and might be preferred for certain indices.

C1307: **Safe haven properties of gold: An application of a dynamic score-driven rotation model**

Presenter: **Ruven Zapf**, Georg-August-Universität Göttingen, Germany

Co-authors: Helmut Herwartz

Gold's safe-haven properties are evaluated through a tail-risk perspective rooted in higher-order moment dynamics. A framework is developed that integrates a multivariate GARCH specification with a dynamic score-driven rotation, enabling the extraction of independent components and time-varying conditional skewness and kurtosis. These higher-order dynamics are translated into one-step-ahead value at risk and expected shortfall forecasts via the Cornish-Fisher expansion, thereby linking distributional asymmetries directly to risk measures. Using daily data for global equities and gold from 2007 to 2025, three main findings emerge. Dynamic rotations substantially improve the calibration of value at risk and expected shortfall forecasts relative to Gaussian, Student-t, and unrotated benchmarks, as verified by standard backtests. Optimal portfolio weights derived from higher-moment risk criteria tilt sharply toward gold during episodes of financial turmoil, including the global financial crisis, COVID-19, and the Russian invasion of Ukraine, while reverting toward equities in tranquil periods. By jointly modeling volatility, skewness, and kurtosis, the analysis reveals a contingent rather than universal safe-haven role for gold, emphasizing that its protective capacity is crisis-specific and state-dependent. The results highlight the economic relevance of higher-order moment dynamics for risk management, portfolio allocation, and the assessment of safe-haven assets.

C1472: **Probabilistic AI for improved tail risk estimation**

Presenter: **Morten Risdal**, NTNU, Norway

Co-authors: Rickard Sandberg

Probabilistic AI models are explored for day-ahead forecasting of return distributions and tail risk. Traditional econometric models rely on restrictive parametric assumptions, while standard machine learning offers only point forecasts without uncertainty. Probabilistic AI addresses these gaps by directly generating flexible conditional return distributions with uncertainty quantification. LSTM and transformer architectures are implemented within a mixture density network (MDN) framework, using realized variance from the CaPiRe dataset and implied volatility from Bloomberg as predictors. The models yield non-parametric return distributions for DJIA constituents, from which value-at-risk (VaR) and expected shortfall (ES) are derived. Performance is compared against econometric and machine learning benchmarks using statistical adequacy tests and scoring metrics. Predictive uncertainty is further decomposed into aleatoric (market risk) and epistemic (model risk) components, enhancing interpretability. Empirical results show that LSTM-MDN and Transformer-MDN consistently outperform benchmarks in distributional accuracy, calibration, and tail risk forecasting, remaining robust across market regimes, including COVID-19. Models incorporating implied volatility deliver the strongest performance, confirming their superior predictive power. Probabilistic AI thus offers a powerful alternative for risk management, portfolio allocation, and derivative pricing.

CC427 Room BCB 307 MACRO-FINANCIAL MODELING

Chair: Demetris Koursaros

C1314: **Simulating macro-financial scenarios of stress for multiple countries in the euro area**

Presenter: **Katrin Assenmacher**, European Central Bank, Germany

Co-authors: Marcel Brautigam, Juan Manuel Figueres, Carla Giglio, Alberto Grassi, Costanza Rodriguez

The purpose is to design an econometric framework to simulate macro-financial scenarios of stress for the four main economies of the euro area. In doing so, a flexible model framework consisting of two blocks is proposed. First, financial shocks are generated via a non-parametric copula estimated on a large dataset of daily financial indicators, which allows generating an index of financial conditions that is used to simulate financial stress events. Second, the joint dynamics of macroeconomic indicators conditional on the copula-based financial shocks in a large multi-country Bayesian VAR model are simulated. Preliminary results show that the proposed framework allows replicating different macro-financial stress scenarios where adverse shocks generated in the financial sector propagate into the overall economy, triggering significant fluctuation in key macroeconomic indicators.

C1165: **Severe weather and financial (in)stability**

Presenter: **Paolo Gelain**, Federal Reserve Bank of Cleveland, United States

Co-authors: Marco Lorusso, M. Marcellino, Claudia Foroni

The effect of severe weather shocks is quantified on the US economy in an environment in which the economy can switch between periods of financial stability and financial instability, like the Great Recession. A New Keynesian dynamic stochastic general equilibrium model with banks and severe weather events is estimated. It is shown that severe weather shocks: 1) have a negative impact on real and financial US variables, sizable only in periods of financial instability, but muted effects on nominal variables; 2) are never a relevant source of business cycle fluctuations; 3) transmit mainly via a deterioration of the quality of capital.

C0199: **A model of sales and promotions**

Presenter: **Demetris Koursaros**, Cyprus University of Technology, Cyprus

The effect of sales and promotions on the pricing decisions of firms is investigated by providing a novel theoretical model where firms face price adjustment costs and offer items on sale to attract bargain hunters. It is demonstrated that in a recession, the frequency of items on sale increases, as well as the effort of consumers to locate such offers. Since the decrease in the cost of living following a recession comes both from price decreases and the combination of more frequent sales and more active bargain hunting by consumers, a price index that simply focuses on prices and neglects high-frequency sales and their weight in the consumers basket appears to be less responsive to shocks. This can explain the low response of inflation in the post-crisis period and thus the breakdown of the Phillips curve as sales items are under-represented in common price indexes. It is also shown that a price index created by placing more weight on sale items in the UK CPI correlates better with the output gap, confirming the model's predictions. Additionally, if agents form inflation expectations using indices that neglect sales and promotions, recessions are exacerbated.

C1533: **On unspanned latent risks in dynamic term structure models**

Presenter: **Tomasz Dubiel-Teleszynski**, University of Liechtenstein, Liechtenstein

Co-authors: Konstantinos Kalogeropoulos, Nikolaos Karouzakis

A parsimonious class of arbitrage-free yields-only Dynamic Term Structure Models (DTSMs) with unspanned latent risks is proposed. We develop an efficient Sequential Monte Carlo (SMC) inferential and prediction scheme that guarantees joint identification of parameters and latent states and takes into account all relevant uncertainties. We use the developed setup to explore the out-of-sample statistical and economic evidence of bond return predictability from the perspective of a real-time Bayesian investor seeking to forecast excess bond returns and maximise her utility. We find that latent factors contain significant predictive power above and beyond the yield curve, offering significant improvement to the out-of-sample predictive performance of models. Most importantly, they exploit information hidden from the yield curve and generate significant utility gains, out-of-sample. Macro-financial linkages are also explored. The hidden component associated with slope risk is countercyclical and links with real

activity.

CC375 Room BCB 310 FUNCTIONAL DATA ANALYSIS
Chair: Matteo Fontana
C0318: Joint analysis of amplitude and phase variations in function-valued traits in quantitative genetics

Presenter: **Yilin Chen**, Kings College London, United Kingdom

Co-authors: Davide Pigoli

Quantitative genetics is the study of continuous traits and the statistical analysis of the relative contributions of genetic and environmental effects to phenotypic variations. Function-valued traits, such as growth curves, are described as a function of a continuous index. Such traits, by nature, can be assessed for continuous genetic variation using a quantitative approach. However, growth curves from different individuals often present variability in both amplitude and phase. A conventional framework can be used, where curves are first registered and then the aligned data are analyzed. Nonetheless, this has the downside of ignoring phase variations in the genetics analysis. To include phase variations in analysis, the proposed method considers the aligned curve and the warping function (mapping the aligned curve to the original observation) as a joint functional object. A functional mixed-effects model that uses functional combined principal components as the data-driven basis is extended to genetically correlated data for this purpose. The proposed procedure is explored in simulations and applied to a real-case study of growth trajectories of red flour beetles.

C1276: Effects of recent and current income on life satisfaction: A functional regression model with endogeneity

Presenter: **Mirkan Oecal**, University of Cologne, Germany

Co-authors: Sven Otto, Dominik Wied

The aim is to propose a functional regression model to analyze both concurrent and historical effects of income on life satisfaction. A model is developed with a recent history operator to capture how past income affects current happiness. A finite-dimensional factor structure is imposed on the functional regressor, reducing the infinite-dimensional problem to a finite-dimensional regression and enabling valid statistical inference. The key contribution is addressing the endogeneity in the income-happiness relationship - particularly reverse causality where happier people may earn more - through an instrumental variables strategy adapted to the factor-augmented model. The methodology builds on recent work in related functional linear models, but extends these by implementing an IV approach for causal interpretation. Estimation proceeds by identifying factors as (left-)singular functions of the cross-covariance operator between life satisfaction and income history, followed by two-stage least squares regression. An empirical application using the GSOEP data demonstrates the practical implementation of the proposed methodology.

C1421: Forecast evaluation with functional data

Presenter: **Shixuan Wang**, University of Reading, United Kingdom

Co-authors: Jack Fosten, Piotr Kokoszka, Tim Kutta

The aim is to propose methods for comparing the accuracy of two competing sets of functional forecasts. This is increasingly relevant as many economic and financial forecasters are interested in predicting variables which are observed as functional data objects. However, to date, there have been no formal statistical tests to evaluate the relative accuracy of competing functional forecasts. A suite of novel tests is proposed, building on the classic Diebold-Mariano test, to provide formal statistical guidance on forecast accuracy in the case of functional data. The asymptotic properties of the tests are derived, as well as self-normalized versions, and the validity of analytic or bootstrap-based critical values is demonstrated. The finite sample performance of the tests is investigated using Monte Carlo simulations. The practical usefulness of the tests is demonstrated by evaluating competing forecasts of U.S. yield curves based on forward rates and a naive benchmark.

C1366: Generalized distance covariance for linearity and independence in functional regression with missing at random responses

Presenter: **Pedro Galeano**, Universidad Carlos III de Madrid, Spain

Co-authors: Manuel Febrero-Bande, Wenceslao Gonzalez-Manteiga

A testing procedure is proposed to jointly assess the linearity between a scalar response and a functional covariate, and the independence between the covariate and the error term, in scalar-on-function regression models with responses missing at random (MAR). The test statistic corresponds to the generalized distance covariance between the functional covariate and the residuals obtained from a linear model fit. To address MAR responses, two slope estimation approaches based on functional principal components (FPCs) are considered: the simplified method, which excludes observations with missing responses, potentially leading to information loss; and (ii) the imputed method, which incorporates additional data by imputing missing responses using the simplified slope estimate. Cross-validation is employed to determine the optimal number of FPCs for each method. The distribution of the test statistic under the null hypothesis is calibrated using residual bootstrap. Monte Carlo simulations show that the proposed procedure can be quite powerful with appropriate choices of the semimetric and its associated parameters for the generalized distance covariance. Additionally, using the imputation method for estimation slightly increases the power of the procedure. The proposed methodology is illustrated through an application involving the modeling of average daily temperatures based on the average number of sunny days recorded at Spanish meteorological stations.

CC432 Room BCB 311 FISCAL AND FINANCIAL ECONOMETRICS
Chair: Mayetri Gupta
C0398: Real time monitoring of global financial stability

Presenter: **Cindy Shin Hwei Wang**, HSBC Business School, Peking University, China

Co-authors: Cheng Hsiao, Yimeng Xie

The aim is to propose an easy-to-implement, online cumulative sum of squares (CUSQ)-type test for monitoring the stability of the global financial system via autoregressive approximation of a mixed panel. The limiting distribution of the test statistics follows a Brownian bridge and is free from model parameters, the bootstrap procedure, or the exact time series properties of each series. Monte Carlo simulations based on various data-generating processes confirm its promising performance. The empirical illustration based on the extension of a prior event-study approach also shows that the likely dates of the break appear to conform with the timing of the occurrence of common events.

C0643: Fiscal impulse responses in a high-dimensional setting

Presenter: **Davide Bucci**, University of Surrey, United Kingdom

The aim is to investigate the relationship between fiscal policy and its macroeconomic effects, focusing on the long-standing debate over the size and sign of the fiscal multiplier. Conventional VAR models struggle in this context because they cannot handle many variables without loss of precision. To address this limitation, the standard specification is replaced with an adaptive LASSO-VAR, which applies data-driven, variable-specific penalties and yields a more flexible yet parsimonious model. A Monte Carlo simulation shows that the adaptive LASSO-VAR recovers impulse-response functions more accurately than the traditional VAR and other leading approaches in the literature, including factor-augmented VARs and Bayesian VARs. In addition, a recoverability test is applied, showing that the proposed model better identifies government-spending shocks. Finally, preliminary empirical results for the United States indicate that the adaptive LASSO-VAR estimates a larger fiscal multiplier than competing models, suggesting stronger real effects of government spending than previously documented. Overall, the evidence confirms that combining adaptive machine-learning techniques with VAR analysis improves identification, inference, and policy-relevant measurement in high-dimensional settings.

C1247: Effects of a money-financed fiscal stimulus without irredeemability of money*Presenter:* **Eiji Okano**, Nagoya City University, Japan*Co-authors:* Masataka Eguchi, Roberto Billi

The purpose is to analyze the effectiveness of a money-financed fiscal stimulus (MF) without irredeemability of money (IM). Although the effectiveness of the MF fiscal stimulus without the IM is weaker than that of the MF fiscal stimulus with the IM, that of the MF fiscal stimulus without the IM is more substantial than a conventional debt-financed (DF) fiscal stimulus. This finding is applicable either in normal times or in a liquidity trap. It assumes not only a closed economy but also a two-country economy. It is found that as the size of a home country increases, the effectiveness of the MF fiscal stimulus without the IM increases, although that of the MF fiscal stimulus with the IM decreases as the size of a home country increases. In addition, it is found that the effectiveness of global MF fiscal stimulus without the IM amid a liquidity trap is more substantial than that of DF fiscal stimulus.

CC422 Room BCB 312 ELECTRICITY MARKETS**Chair: Robinson Kruse-Becher****C1237: Forecasting electricity prices using bid data in times of distress***Presenter:* **Ainhoa Zarraga**, University of the Basque Country, Spain*Co-authors:* Aitor Ciarreta, Blanca Martinez

In the context of ongoing market reforms within the European Union, developing more accurate electricity price forecasting techniques is becoming increasingly essential for all market participants. The increasing penetration of renewables, coupled with the geopolitical tensions between exporters and importers of raw energy materials that emerged in 2021, has resulted in growing pressure on investors to adopt effective hedging strategies to protect their assets in this turbulent economic environment. The aim is to propose a price forecasting approach based on publicly available auction data to fit the supply and demand electricity curves of the Iberian electricity market for the period 2021-2024. First, fractional polynomial and logistic functions are fit to historical sales and purchase bidding data to estimate the equilibrium prices. Secondly, several time series models are specified and estimated for the historical and estimated prices. Thirdly, a rolling window is used to estimate models for both time series prices and forecast one-day-ahead prices for 2024. In-sample and out-of-sample error criteria are used. Results show that using fractional polynomials and logistic functions accurately replicates the observed prices. The out-of-sample forecasting analysis shows that the fractional polynomial functions outperform not only the naive model, but also the models based on historical prices.

C1289: Electricity price forecasting in the day-ahead market: Averaging forecasts vs breakpoint detection*Presenter:* **Piotr Zaborowski**, Wroclaw University of Science and Technology, Poland*Co-authors:* Rafal Weron

The profitability of a battery energy storage system (BESS) in the day-ahead market is determined by the precise timing of buying (charging) and selling (discharging) electricity. The latter requires reliable electricity price forecasts for the next day. Two approaches are compared using regression-based models. The first averages forecasts across calibration windows of varying lengths, balancing short-term responsiveness with long-term stability. The second uses structural break detection with the pruned exact linear time (PELT) algorithm, calibrating models only within homogeneous segments. The approaches are evaluated using both statistical (MAE, RMSE) and economic criteria (average opportunity cost of a BESS trading strategy and the Sharpe ratio) on four European day-ahead markets: Germany, Finland, Poland, and Spain. The test period includes the COVID-19 pandemic and the 2022 energy crisis. The results show that averaging-based models outperform breakpoint-based approaches in both accuracy and profitability.

C1331: Stealing accuracy: Predicting day-ahead electricity prices with temporal hierarchy forecasting (THieF)*Presenter:* **Rafal Weron**, Wroclaw University of Science and Technology, Poland*Co-authors:* Arkadiusz Lipiecki, Kaja Bilinska, Nikolaos Kourentzes

The purpose is to introduce the concept of temporal hierarchy forecasting (THieF) in predicting day-ahead electricity prices and show that reconciling forecasts for hourly products, 2- to 12-hour blocks, and baseload contracts significantly (up to 13%) improves accuracy at all levels. These results remain consistent throughout a challenging 4-year test period (2021-2024) in the German power market and across model architectures, including linear regression, a shallow neural network, gradient boosting, and a state-of-the-art transformer. Given that (i) trading of block products is becoming more common and (ii) the computational cost of reconciliation is comparable to that of predicting hourly prices alone, it is recommended to use it in daily forecasting practice.

C1443: European electricity market integration: A volatility spillover approach*Presenter:* **Cristina Pizarro-Irizar**, University of the Basque Country, Spain*Co-authors:* Aitor Ciarreta, Evelyn Lizeth Chanatasig, Ainhoa Zarraga

The purpose is to apply a forecast-error variance decomposition approach, including both variances and covariances, for analyzing price volatility spillover effects across nine European electricity markets from 2015 to 2023. Static and dynamic analyses of aggregated spillover effects and their directional decomposition between markets are conducted. The static analysis shows a high degree of volatility transmissions, which indicates a high level of market integration, though there are differences from one country to another. The dynamic analysis shows that the evolution of the total volatility spillovers can be mainly related to the deployment of renewable electricity penetration, the increase in electricity flows between markets, and the exceptional increase in the natural gas price following the economic recovery after the pandemic, exacerbated by the war in Ukraine. Findings provide useful information to regulators to improve current and future policies regarding the European Union's goals for market integration and renewable energy penetration.

CC388 Room BCB 403 CATEGORICAL DATA**Chair: Alejandro Murua****C1180: Ordinal item response theory models for the evaluation of heterogeneity in attitudes***Presenter:* **Ingrid Maurer**, University of Malaga, Spain*Co-authors:* Gerhard Tutz

The contribution provides ordinal item response models to investigate heterogeneity in attitudes. In contrast to standard models, the approach accounts for covariate effects in studying two essential components of attitudes: direction and strength. An application to the evaluation of presidential candidate traits exemplifies the types of substantive insights that can be gained.

C1299: Quantifying and testing dependence to categorical variables*Presenter:* **Daniel Strenger-Galvis**, Graz University of Technology, Austria*Co-authors:* Siegfried Hoermann

The purpose is to suggest a dependence coefficient between a categorical variable and some general variable taking values in a metric space. Important theoretical properties are derived, and the large sample behavior of the suggested estimator is studied. Moreover, an independence test is developed which has an asymptotic χ^2 -distribution if the variables are independent, and it is proven that this test is consistent against any violation of independence. Some extensions are discussed, including a variant of the coefficient for measuring conditional dependence.

C1321: Aitchison geometric characterization of quasi-symmetry in square contingency tables*Presenter:* **Keita Nakamura**, Tokyo University of Science, Japan

Co-authors: Tomoyuki Nakagawa, Kouji Tahata

Quasi-symmetry provides a flexible alternative to complete symmetry in the analysis of square contingency tables. Within Aitchison geometry, tables are treated as compositions and analyzed through the centered log-ratio transform, which endows the simplex with a Euclidean structure suitable for distances, norms, and projections. A table is quasi-symmetric when its interaction component coincides with its transpose, and the set of quasi-symmetric tables forms a linear subspace. An explicit orthogonal projection maps any table to its nearest quasi-symmetric approximation by averaging the interaction component with its transpose while preserving the row and column geometric marginals. Departures from quasi-symmetry are summarized by the simplicial quasi-skewness, defined as the squared Aitchison norm of the skew-symmetric interaction table that is orthogonal to the quasi-symmetric subspace. The quasi-skewness array highlights the signed cell-wise contributions as relative proportions of the total quasi-skewness, thereby quantifying the relative importance of each departure. An application to real data will be presented to illustrate the practical use of the proposed approach.

C1341: Information-geometric viewpoint on partitioning of test statistics for symmetry in contingency tables

Presenter: **Tomoyuki Nakagawa**, Meisei University & RIKEN Center for Brain Science, Japan

Co-authors: Takeru Matsuda, Kouji Tahata

Numerical and asymptotic partitioning of goodness-of-fit statistics has been investigated for numerous models in contingency tables. The purpose is to explore how goodness-of-fit statistics for symmetry in contingency tables can be partitioned. The symmetry model is dually flat and can be characterized as the intersection of an e-flat submodel and an orthogonal m-flat submodel. Based on this result, the Wald statistic for the symmetry model can be exactly partitioned into components corresponding to the submodels. On the other hand, there are very few models for which the likelihood ratio test statistic for the symmetry model can be exactly partitioned. The focus is on the relationship between information geometry and the exact partitioning of goodness-of-fit statistics for symmetry in contingency tables.

CC387 Room BCB 407 CAUSAL INFERENCE AND NETWORK DATA

Chair: Steven Gilmour

C1304: Full of noise? A bootstrap confidence interval that controls size in Diebold-Yilmaz networks

Presenter: **Milan Csaba Badics**, Corvinus University of Budapest, Hungary

Over the past decade, Diebold-Yilmaz (DY) connectedness has become a workhorse for empirical network analysis, yet inference has received relatively little attention. The few papers that attempt it typically adopt Efron's percentile bootstrap confidence intervals but rarely evaluate whether those tests achieve nominal size, despite its importance for policy and trading. It is demonstrated that even under favorable conditions, long samples, low persistence, homoskedastic innovations, and unbiased impulse responses, the DY connectedness measures exhibit skewness and bias that inflate rejection rates. A Hall-type bootstrap confidence interval is proposed with coverage calibration via bootstrap-after-bootstrap, delivering size-correct inference for the connectedness table. In Monte Carlo simulations, the proposed interval achieves an empirical size close to the nominal level. The method reduces spurious link detections, supporting more credible policy analysis and more reliable minimum-connectedness portfolio construction.

C1470: Inferring age-stratified social networks from contact data with Gaussian mixtures

Presenter: **Luke Murray Kearney**, University of Warwick, United Kingdom

Co-authors: Emma Davis, Matt J Keeling

Capturing the structure of a population and characterising contacts within the population are key to reliable projections of infectious diseases. Two main elements of population structure, contact heterogeneity and age, have been repeatedly demonstrated to be key in infection dynamics, yet are rarely combined. While there are a few key examples of contact networks being measured explicitly, in general, there is a need to construct the appropriate surrogate networks from individual-level data. Using data from open-source social contact surveys, an algorithm is developed to generate an extrapolated network that preserves both age-structured mixing and heterogeneity in the number of contacts. The spread of infection is then simulated through the population, constrained to have a given basic reproduction number, R_0 . Given the over-dominant role that highly connected nodes ('superspreaders') would otherwise play in early dynamics, transmission is scaled by the duration of contacts, providing a better match to surveillance data for numbers of secondary cases. Showing that for COVID-like parameters, including both heterogeneity and age-structure, reduces epidemic size. A robust methodology, therefore, allows for the inclusion of the full wealth of data commonly collected by surveys but frequently overlooked to be incorporated into more realistic transmission models of infectious diseases.

C1283: Estimating effects of time-varying continuous-treatment in regression discontinuity designs

Presenter: **Shoichiro Nagasaka**, Doshisha University, Japan

Co-authors: Yutaka Kano

The purpose is to study statistical causal inference when treatment assignment is made based on a threshold, treatment intensity is continuous, and effects evolve over time. Building on regression discontinuity designs (RDD), it is noted that continuous-treatment approaches, including the framework of a prior study, identify effects at the cutoff but do not explicitly model temporal dynamics. A time-varying continuous-treatment RDD is introduced that preserves identification from the running variable crossing a known cutoff yet allows treatment to vary across time. Estimation combines kernel-weighted local methods tailored to dynamic settings, local quantile regression to recover distributional impacts that may shift, and local linear regression to estimate average responses with reduced boundary bias near the threshold. Together, these estimators accommodate continuous dosage, exploit local smoothness, and trace time-dependent treatment response patterns. Assumptions for identification and consistency are stated, bandwidth selection is discussed, and it is shown how the approach nests the standard RDD as a special case. Monte Carlo studies and empirical evidence demonstrate practical relevance, including sensitivity to bandwidth and kernel choices. Code and replication materials are released. Applications include clinical and economic panels with evolving policies and thresholds measured repeatedly over time across units and periods.

C1209: Causal DAG identification for count data via Poisson-Thinning structural equation models

Presenter: **Penggang Gao**, Kyoto University, China

Co-authors: Ming Cai, Hisayuki Hara

Count data arise frequently in epidemiology, economics, and social sciences, motivating the development of causal models that account for discrete distributions. Existing Poisson branching approaches encode causal structure via binomial thinning operators but face two major limitations: The presence of equivalence classes that hinder full identifiability, and the restriction of exogenous disturbances only to the Poisson distribution. To address these issues, a structural equation model based on the Poisson thinning operator is introduced. The model is proven fully identifiable under exogenous noise following Poisson, negative binomial, or zero-inflated Poisson distributions. The framework includes a parameter estimation procedure derived from the structural equations and a likelihood score-based identification procedure for orienting edges in candidate graphs. Extensive simulations with exogenous noise from these distributions show high identification precision in recovering causal structures, and high accuracy in parameter estimation under finite samples. In summary, the proposed framework advances models and identification theory for count data, providing practical methodology for causal discovery in observational studies.

CV434 Room Virtual R01 STATISTICAL MODELS AND INFERENCE

Chair: Ralf Wilke

C1342: A computational note on the penalized correlation-based estimators in linear and generalized linear models

Presenter: **Mina Norouzirad**, Center for Mathematics and Applications (NOVA Math), NOVA School of Science and Technology, Portugal

Co-authors: Marta Lopes, Tomas Bandeira, Ricardo Moura

Penalized regression methods such as ridge, lasso, and elastic net have become standard tools in statistical modeling, yet they do not explicitly account for correlations among explanatory variables. This limitation may reduce performance in the presence of strong multicollinearity, leading to unstable or biased coefficient estimates. To address this issue, a prior study proposed a correlation-based penalty that incorporates dependence structures directly into the estimation process. Despite its theoretical appeal, the practical implementation of this estimator is challenging, particularly in high-dimensional settings where the number of predictors can exceed the sample size. The purpose is to revisit correlation-based penalized regression from a computational perspective. The aim is to investigate efficient algorithms for estimating this estimator, with particular attention to numerical stability, and to compare their computational time and performance across different scenarios. The focus is on exploring approaches that may render the estimator feasible and useful in practice, especially in high-dimensional contexts. The ultimate objective is to provide a practical implementation, thereby extending the toolbox of penalized regression methods available for linear and generalized linear models.

C0996: Modeling correlated ordinal data through Students t copula

Presenter: **Alessandro Barbiero**, Università degli Studi di Milano, Italy

The t copula is a member of the family of elliptical copulas, which includes, among others, the Gaussian copula, with which it shares many properties. Unlike the latter, which is tail-independent, the t copula always exhibits equal and non-null lower and upper tail-dependence coefficients, even when its components are uncorrelated. Although the use of copulas is less straightforward in the discrete case, primarily due to theoretical challenges related to Sklar's theorem, several contributions in the statistical literature have explored the use of the Gaussian copula for constructing multivariate discrete distributions. The use of the t copula is illustrated for modeling bivariate correlated discrete data by investigating two aspects. First, it is examined how its additional parameter, the degrees of freedom, can lead to different bivariate discrete distributions when the correlation parameter is held fixed. To quantify the differences between distributions, suitable distance measures are employed. Second, for a fixed (positive) correlation, a fixed number of degrees of freedom, and a fixed cardinality k of the support of the two discrete distributions, it is numerically identified which common k -point probability distribution induces the maximum linear correlation. It is found that the resulting correlation between the discrete random variables can exceed the copula's correlation parameter: A phenomenon that contrasts with the behavior observed when using the Gaussian copula.

C1105: Comparative study on tail probability estimation in i.i.d. settings

Presenter: **Taku Moriyama**, Yokohama City University, Japan

Tail probability estimators in i.i.d. settings are considered. There are mainly two ways for the estimation: The fitting to the generalized Pareto distribution and the fully nonparametric estimation. The fitting estimator is justified by the approximation proven in the extreme value theory; however, the accuracy depends on the target point, i.e., how extremely large the target is. The nonparametric estimator does not need the approximation and has the advantage of wide applicability; however, the optimal regularization parameter depends on both the target point and the extreme value index. Both theoretical and numerical comparative studies on tail probability estimation are conducted. Asymptotic convergence rates of estimators are obtained, and the mean integrated squared errors are numerically surveyed by a simulation study.

C1505: Goodness-of-fit tests for cure rate quantile regression

Presenter: **Mercedes Conde-Amboage**, University of Santiago de Compostela, Spain

Co-authors: Wenceslao Gonzalez-Manteiga, Cesar A Sanchez-Sellero

In classical survival analysis, a fundamental assumption is that all individuals will eventually experience the event of interest. However, it often occurs that a subset of subjects will never experience the event. These individuals are typically considered to have infinite survival times and are classified as 'cured'. To deal with this phenomenon, classical survival models have been extended to what is commonly referred to as cure models. Throughout this talk, a new lack-of-fit test for cure models in the context of quantile regression is presented. This new proposal represents the first contribution in the literature to test the effect of a set of covariates on survival time using empirical processes based on residuals. The asymptotic behaviour of the test statistics will be derived. In addition, an extensive simulation study and a real-data application will be presented to show the performance of the new proposal in practice.

Monday 15.12.2025

14:40 - 16:20

Parallel Session N – CFE-CMStatistics 2025

CI006 Room BCB 206 CFE SPECIAL INVITED SESSION: A TRIBUTE TO H. PESARAN I**Chair: Robert Taylor****C0160: The effects of macroeconomic shocks on firms' price and wage inflation expectations***Presenter:* **Martin Weale**, Kings College London, United Kingdom*Co-authors:* Tomasz Wieladek

The purpose is to study the effects of macroeconomic shocks on firm-level price and wage inflation expectations in the UK's CBI Industrial Trends Survey. A novel feature is that reduced-form forecast errors are mapped into structural supply and demand shocks for the analysis. The dataset has information on wage, alongside price, inflation expectations from 2012Q2-2024Q3. This allows exploring the differences in wage and price inflation expectation formation before the pandemic and to examine whether the influence of shocks was different from 2020Q2 onwards, from what had been experienced earlier. It is found that supply shocks tend to have more influence on wage and price expectations than do demand shocks. There were marked changes in the effects of both types of shocks during the period from 2020Q2 onwards.

C0154: Model-based test for asset price bubbles*Presenter:* **Robert Taylor**, University of Essex, United Kingdom

The purpose is to develop tests for asset price bubbles that are derived directly from the standard stock pricing equation commonly employed in the finance literature, whereby stock prices are determined by expectations of the future price, future dividends, and an unobserved component. The general solution to this equation shows that the price of an asset is the sum of a fundamental component and an explosive component that only takes non-zero values during bubble periods. While the current literature focuses almost exclusively on modeling asset prices as a single autoregressive process and testing the null of a unit root against the alternative of explosivity, this is not the model implied by finance theory. Instead, the solution is directly used to the asset price equation, and locally best invariant [LBI] motivated tests is constructed, designed to test the null that the innovations to the bubble component of the asset pricing equation have zero variance (and hence no bubble exists), against the alternative that these innovations have non-zero variance (and hence an asset price bubble exists). Recursive versions of the proposed test are also developed, and it is shown via simulation that these tests are significantly more powerful than the industry standard tests developed by a prior study.

C1007: Jackknife inference with two-way clustering*Presenter:* **James MacKinnon**, Queen's University, Canada*Co-authors:* Morten Nielsen, Matt Webb

For linear regression models with cross-section or panel data, it is natural to assume that the disturbances are clustered in two dimensions. However, the finite-sample properties of two-way cluster-robust tests and confidence intervals are often poor. Several ways to improve inference with two-way clustering are discussed. Two of these are existing methods for avoiding, or at least ameliorating, the problem of undefined standard errors when a cluster-robust variance matrix estimator (CRVE) is not positive definite. One is a new method that always avoids the problem. More importantly, a family of new two-way CRVEs is proposed based on the cluster jackknife, and it is proven that they yield valid inferences asymptotically. Simulations for models with two-way fixed effects suggest that, in many cases, the cluster-jackknife CRVE combined with the new method yields surprisingly accurate inferences. A simple software package is provided, *twowayjack* for Stata, that implements the recommended variance estimator.

CI016 Room BCB 307 CAUSAL INFERENCE UNDER UNOBSERVED CONFOUNDING**Chair: Mark Steel****C0619: Challenges in causality: Sensitivity analysis and automation***Presenter:* **Carlos Cinelli**, University of Washington, United States

Causal inference has emerged as a central research area at the intersection of statistics and computer science, with numerous applications in empirical fields such as biomedicine and the social sciences. However, despite significant progress over the past three decades, many important problems remain open. The aim is to highlight two critical research directions: (i) drawing valid conclusions when key assumptions, such as the absence of unobserved confounding, are violated, and (ii) automating the causal inference pipeline. Key open challenges are outlined in the areas of sensitivity analysis and partial identification, as well as the automation of such tasks. Addressing these challenges is essential for making robust causal inferences under realistic settings.

C0695: Bayesian model averaging in causal instrumental variable models*Presenter:* **Gregor Steiner**, University of Warwick, United Kingdom*Co-authors:* Mark Steel

Instrumental variables are a popular tool to infer causal effects under unobserved confounding, but choosing suitable instruments is challenging in practice. The aim is to propose gIVBMA, a Bayesian model averaging procedure that addresses this challenge by averaging across different sets of instrumental variables and covariates in a structural equation model. The approach extends previous work through a scale-invariant prior structure and accommodates non-Gaussian outcomes and treatments, offering greater flexibility than existing methods. The computational strategy uses conditional Bayes factors to update models separately for the outcome and treatments. It is proven that this model selection procedure is consistent. By explicitly accounting for model uncertainty, gIVBMA allows instruments and covariates to switch roles and provides robustness against invalid instruments. In simulation experiments, gIVBMA outperforms current state-of-the-art methods. Its usefulness is demonstrated in an empirical application, estimating the causal effect of education on income.

C1408: Potential outcome modeling and estimation in DiD designs with staggered treatments*Presenter:* **Siddhartha Chib**, Washington University in Saint Louis, United States*Co-authors:* Kenichi Shimizu

The aim is to propose the first potential outcome modeling of difference-in-differences designs with multiple time periods and variation in treatment timing. Importantly, the modeling respects the two key identifying assumptions: parallel trends and no-anticipation. A straightforward Bayesian approach is then introduced for estimation and inference of the time-varying group-specific Average Treatment Effects on the Treated (ATT). To improve parsimony and guide prior elicitation, the model is reparametrized in a way that reduces the effective number of parameters. Prior information about the ATT is incorporated through black-box training sample priors and, in small-sample settings, by thick-tailed t-priors that shrink ATT of small magnitudes toward zero. A computationally efficient Bayesian estimation procedure is provided, and a Bernstein-von Mises-type result is established that justifies posterior inference for the treatment effects. Simulation studies confirm that the method performs well in both large and small samples, offering credible uncertainty quantification even in settings that challenge standard estimators. The practical value of the method is illustrated through an empirical application that examines the effect of minimum wage increases on teen employment in the United States.

CO081 Room BCB G07 HiTEC: CLUSTERING OF COMPLEX DATA STRUCTURES**Chair: Maria Brigida Ferraro****C0504: Modeling residual heterogeneity within the mixture of experts framework***Presenter:* **Francesca Martella**, La Sapienza University of Rome, Italy*Co-authors:* Dalila Failli, Maria Francesca Marino

The mixture of latent trait analyzers (MLTA) is a model-based clustering framework tailored for multivariate categorical data, combining features from latent class and latent trait models. It allows for both the clustering of units and the modeling of residual within-cluster variability through continuous latent variables. A novel extension of the MLTA model that incorporates the effects of covariates (concomitant variables) is proposed in a flexible and comprehensive way. Specifically, covariates are allowed to influence cluster formation, the distribution of the response variables, both, or neither, mirroring the structure of standard mixtures of expert models. This generalization enhances the interpretability of clustering results and captures more complex data structures, especially in contexts where observed characteristics are expected to affect both group membership and response behavior. The model is estimated via an EM-type algorithm using variational approximations. A motivating example from the biomedical domain illustrates how the proposed model can simultaneously uncover latent patient profiles and account for the effects of clinical covariates on both clustering and outcomes. The flexibility of this approach makes it suitable for a wide range of applications in social sciences, health, and beyond.

C0657: Model-based clustering of population of networks via extended stochastic block models

Presenter: **Maria Francesca Marino**, University of Florence, Italy

Co-authors: Monia Lupporelli, Giulia Capitoli

The population of networks arises when interactions between nodes are observed repeatedly over a set of units. The stochastic block model is extended to provide a joint clustering of units and network nodes. This is done by considering in the model specification a set of unit- and node-specific, discrete, latent variables, able to capture dependences among the observed data. A former set of latent variables allows partition units into classes based on the observed, unit-specific, network features. A latter set of latent variables is employed to identify clusters of stochastically equivalent nodes sharing similar connectivity profiles. Parameter estimation is conducted within a maximum likelihood framework. However, deriving the likelihood function is a challenging task, as it would require the solution of a multiple summation over all latent variables in the model. To solve the issue, the use of appropriate approximation methods is required. The effectiveness of the proposal is evaluated both via simulations and a real data application from the medical field.

C0700: Robust clustering using maximized mutual information

Presenter: **Mackenzie Neal**, McMaster University, Canada

Co-authors: Paul McNicholas, Arthur White

Flow cytometry technology allows for the analysis of disease impact on individual cells. Population identification in flow cytometry datasets is commonly performed, with manual gating being the most frequently used method. However, manual gating is subjective, often irreproducible, and becoming increasingly more difficult to perform as the complexity of flow cytometry experiments and datasets grows. Methods have been proposed to automate population identification of flow cytometry data; many such methods rely on generative clustering models. Ideas from both generative and discriminative models are incorporated to present a clustering algorithm capable of capturing irregular sub-populations common in biological datasets. Although flow cytometry is the primary motivation for the work proposed herein, this method can be applied to any dataset wherein the goal is to obtain intuitive clustering solutions. This is demonstrated by comparison to popular clustering methods on various simulated and real datasets.

C1141: Parsimonious ultrametric Manly mixture models

Presenter: **Paul McNicholas**, McMaster University, Canada

Co-authors: Alexa Sochaniwsky

A family of parsimonious ultrametric mixture models with the Manly transformation is developed for clustering high-dimensional and asymmetric data. Advances in Gaussian mixture modelling sufficiently handle high-dimensional data but struggle with the presence of cluster skewness. Also, while these advances reduce the number of free parameters, they often provide limited insight into the structure and interpretation of the clusters. To address this shortcoming, the extended ultrametric covariance structure and the Manly transformation are used, resulting in the parsimonious ultrametric Manly mixture model family. Model selection proves challenging, and a two-step model selection procedure is proposed. Simulation studies and real data analyses are used for illustration.

CO182 Room BCB G08 EMPIRICAL MACRO

Chair: Michael Owyang

C0201: An order and scale-invariant identification scheme for structural dynamic models

Presenter: **Neville Francis**, UNC Chapel Hill, United States

The purpose is to introduce OASIS, an order- and scale-invariant identification scheme for structural dynamic models, which includes structural vector autoregressions (SVARs), local projections, and impulse-response function matching within DSGE frameworks. OASIS identifies the orthogonal rotation that maximizes the correlation between structural shocks and their corresponding reduced-form innovations. This objective is similar to that of a Cholesky decomposition, but the latter is constrained by the ordering of variables. Additional theoretical insights regarding Cholesky identification are also provided. To evaluate OASIS against Cholesky, a significant number of SVAR studies are revisited. The correlations between structural and reduced-form shocks are generally very large. Notably, it is found that OASIS effectively resolves the price puzzle observed in impulse responses to monetary policy shocks.

C0224: Evaluating the dynamics of monetary policy: The implications of speed, level, and duration in interest rate adjustments

Presenter: **Jonghyuck Lim**, The University of North Carolina at Chapel Hill, United States

Co-authors: Neville Francis, Michael Owyang

The aim is to examine how oil supply shocks reshape the distribution of U.S. household inflation expectations. Using a functional vector autoregression (fVAR) framework, the evolution of macroeconomic variables and the full cross-sectional density of household expectations are jointly modeled. Identification follows the Bayesian proxy SVAR approach with external oil supply instruments from a prior study. The results reveal that oil shocks disproportionately shift the right tail of the expectations distribution, raising upper-tail beliefs persistently while lower-tail beliefs revert quickly, leading to widening disagreement and long-lived heterogeneity. These findings highlight the limits of representative-expectations frameworks and underscore the importance of monitoring and communicating to the entire distribution of expectations to preserve policy credibility.

C0446: How election shocks impact markets: Evidence from sectoral stock prices

Presenter: **Aaron Amburgey**, Northwestern University, United States

The effects of U.S. presidential election cycles are examined on sectoral stock markets. Using a high-frequency identification approach, a novel "election shock" series are constructed, which captures exogenous surprises in election probabilities. Aside from election outcomes, the largest shocks are associated with events that are orthogonal to innovations in the macroeconomy, e.g., scandals and debates. These shocks have immediate effects on asset prices in sectors that are differentially impacted by the policy platforms of the two major U.S. political parties. In particular, shocks favoring Republican (Democratic) candidates increase (decrease) the asset prices in the energy and defense sectors, while decreasing (increasing) prices in the clean energy sector. These effects persist overtime.

C0585: Housing bubble contagion across US cities

Presenter: **Daniel Soques**, University of North Carolina Wilmington, United States

Co-authors: Michael Owyang, Jeremy Piger

A multi-region qualitative vector autoregression (QualVAR) model is developed to study the transmission of housing bubbles across US cities. Each city transitions between fundamental and bubble regimes based on a continuous latent variable that depends on lagged regime indicators from other cities. The model allows for the possibility that the likelihood of entering a bubble regime in one city is influenced by the regime history of others, enabling the detection of contagion in housing market dynamics. Fundamental house prices are determined by regional economic drivers, while bubbles are defined as persistent deviations from these fundamentals. Bayesian estimation is conducted using panel data on MSA-level housing prices and regional fundamentals. The model provides a tractable structure for analyzing the spatial diffusion of housing bubbles and supports the real-time identification of speculative episodes. This approach contributes to the understanding of interdependence across regional housing markets by allowing for regime interactions in a dynamic, multi-region setting.

CO120 Room BCB G09 RISK ANALYSIS AND MODEL SPECIFICATION
Chair: Jonas Andersson
C0207: Extending modification indices

Presenter: **Johan Lyhagen**, Uppsala University, Sweden

Co-authors: Yushu Li

Thematic theories are incomplete, meaning that they do not give a full description of how to specify a model. Rather, they focus on certain aspects of interest, which means that there are parts of the model that need to be empirically decided. This includes lag-lengths in time series analysis, factor structures, and correlations amongst errors in SEM, or control variables in causal inference (sensitivity analysis). In SEM, there are modification indices, mainly for the purpose of improving the fit of the model, that estimate the increase in the likelihood when relaxing a restriction. Subsequently, one can also derive an estimated parameter change when relaxing a restriction. This is generalized to relaxing more than one parameter, focusing on the estimated parameter change in the parameters of interest, and deriving this in the GLM setting as well as in the traditional SEM. Theoretical results and Monte Carlo simulations are included to investigate the small sample properties, and empirical examples to show the usefulness for empirical researchers.

C0271: A local Gaussian correlation block bootstrap test and the NordLink cables effect on electricity prices

Presenter: **Baard Stoeve**, University of Bergen, Norway

Co-authors: Sondre Hoelleland, Kristian Gundersen, Jan Bulla

Financial returns often exhibit strong autocorrelation and complex nonlinear dependencies, making standard bootstrap methods unsuitable due to their assumption of independence. To address this, prior studies have used local Gaussian correlation (LGC) and GARCH filtering to capture asymmetric and nonlinear dependence structures. However, GARCH filtering can be cumbersome and may distort the data. A block bootstrap approach is proposed that preserves the nonlinear and time-varying dependence in financial data without requiring filtering. A robust test for financial contagion is also introduced, designed to work with this bootstrap method. The method is illustrated, using price data from the German and Norwegian electricity markets, focusing on whether dependence structures changed after the NordLink interconnector was established in May 2021, a context known for nonlinear and contagious price dynamics.

C0388: Using precipitation forecasts to predict insurance claims

Presenter: **Sondre Hoelleland**, Norwegian School of Economics, Norway

Co-authors: Haakon Otneim, Etienne Dunn-Sigouin, Mahsa Gorji, Geir Drage Berentsen

Climate change is affecting insurers worldwide, as more extreme weather leads to increasingly severe and frequent damage to infrastructure. A framework for short-term property insurance claim forecasting is presented, which facilitates early customer warnings and efficient resource allocation. Using ensemble precipitation forecasts, weather-driven claims in Bergen and Oslo are modeled as a rare-event binary classification problem, employing probabilistic regression and machine learning methods. Models are evaluated on discrimination, reliability, and economic value. Results indicate that using weather forecasts enhances model discrimination, reliability, and operational cost efficiency, compared with baseline scenarios. Case studies of extreme rainfall events illustrate practical application, demonstrating how insurers can leverage publicly available forecasts to anticipate increased claim risks and improve their response strategies.

C0518: The risk aversion coefficient in the tangency portfolio for large dimensions and a singular covariance matrix

Presenter: **Stanislas Muhinyuza**, Linnaeus University, Sweden

Co-authors: Peter Karlsson, Stepan Mazur

The finite-sample distribution of the risk aversion coefficient of the tangency portfolio is derived in the form of a stochastic representation (SR). The focus is on the situation where both the population and the sample covariance matrices of asset returns are singular, particularly when the portfolio size is larger than the number of observations and the returns are identically and independently distributed. The derived SR is used to derive the moments of the risk aversion coefficient and to establish its high-dimensional approximation. Furthermore, through simulation, the good performance of the proposed high-dimensional approximation is documented.

CO271 Room Virtual R01 STATISTICAL MODELING AND MACHINE LEARNING FOR COMPLEX DATA
Chair: Dhrubajyoti Ghosh
C0575: Unified Bayesian nonparametric framework for ordinal, survival, and density regression using complementary log-log link

Presenter: **Entejar Alam**, University of Texas at Austin, United States

Co-authors: Antonio Linero

Applications of the complementary log-log (cloglog) link to problems in Bayesian nonparametrics are developed. Although less commonly used than the probit or logit links, it is found that the cloglog link is computationally and theoretically well-suited to several commonly used Bayesian nonparametric methods. The starting point is a Bayesian nonparametric model for ordinal regression. It is first reviewed how the cloglog link uniquely sits at the intersection of the cumulative link and continuation ratio approaches to ordinal regression. Then, a convenient computational method is developed for fitting these ordinal models using Bayesian additive regression trees. Next, the ordinal regression model is used to build a Bayesian nonparametric stick-breaking process and show that, under a proportional hazards assumption, the stick-breaking process can be used to construct a weight-dependent Dirichlet process mixture model. Again, Bayesian additive regression trees lead to convenient computations. These models are then extended to allow for Bayesian nonparametric survival analysis in both discrete and continuous time. The models have desirable theoretical properties, and this is illustrated by analyzing the posterior contraction rate of the ordinal models. Finally, the practical utility of the cloglog models is demonstrated through a series of illustrative examples.

C0577: A novel longitudinal rank-sum test for multiple primary endpoints in clinical trials

Presenter: **Dhrubajyoti Ghosh**, Kennesaw State University, United States

Neurodegenerative disorders such as Alzheimer's disease (AD) present a significant global health challenge, characterized by cognitive decline, functional impairment, and other debilitating effects. Current AD clinical trials often assess multiple longitudinal primary endpoints to evaluate treatment efficacy comprehensively. Traditional methods, however, may fail to capture global treatment effects, require larger sample sizes due to multiplicity adjustments, and may not fully exploit multivariate longitudinal data. To address these limitations, the longitudinal rank sum test (LRST) is introduced, a novel nonparametric rank-based omnibus test statistic. The LRST enables a comprehensive assessment of treatment efficacy across multiple endpoints and time points, without the need for multiplicity adjustments, thereby effectively controlling Type I error while enhancing statistical power. It offers flexibility against various data distributions encountered in AD research and maximizes the utilization of

longitudinal data. Extensive simulations and real-data applications demonstrate the LRST's performance, underscoring its potential as a valuable tool in AD clinical trials.

C0629: **Signal detection under unknown background when only one unlabeled data is available**

Presenter: **Aritra Banerjee**, University of Minnesota, United States

Co-authors: Sara Algeri, Lydia Brenner, Oliver Rieger

Searches for new physics involve detecting the presence of a specific signal in data that is contaminated by a background arising from several other sources. This task is particularly challenging when a reliable description of the background is unavailable. The aim is to develop a statistical method to test the presence of the signal of interest in the data and to estimate the signal proportion even when the background is unknown or misspecified. Moreover, a signal search is proposed only using a single physics dataset generated from the experiments that may or may not contain the signal of interest. The approach relies on using orthonormal expansion to model the deviation between a proposal density and the unknown density generating the data. It is proposed to choose the proposal density in such a way that one can ensure a conservative estimate of the signal proportion to avoid false discovery. Reliability of this approach is demonstrated through simulation studies, application on realistic simulated data from the Fermi Large Area Telescope, and on data from the ATLAS experiment. A comparative analysis is also performed of the proposed method with the spurious signal method commonly employed in particle physics, and cases are explored where the latter leads to false discoveries.

C0763: **Streaming prediction with hash function based methods**

Presenter: **Aleena Chanda**, University of Nebraska - Lincoln, United States

The traditional empirical distribution function (EDF) becomes computationally and memory-intensive for large data streams, limiting their utility in real-time applications. A novel method is proposed for estimating a distribution function in streaming data built on the count-min sketch algorithm. The estimated empirical distribution function (EEDF) overcomes these limitations by using probabilistic hash functions to approximate frequency distributions with limited memory effectively. By dynamically adjusting histogram interval lengths, the method provides fine-grained approximations of the empirical distribution without storing all data points. The algorithm operates in a single pass and maintains computational efficiency, making it well-suited for streaming settings. While the count-min sketch exhibits a slight upward bias, particularly for low-frequency elements, the effect is minimal compared to its scalability advantages. In predictive contexts, the median of the EEDF typically performed better than most other predictors tested but roughly tied with Gaussian process prior (GPP) predictors that included a bias term. It is concluded that Bayesian and Bayes-like methods are typically among the most effective approaches for prediction in M-open settings.

CO252 Room BCB 207 ADVANCES IN FISCAL POLICY

Chair: Davide Furceri

C1006: **AI meets fiscal policy: Fiscal actions across 140+ countries**

Presenter: **Adrian Peralta Alva**, IMF, United States

Co-authors: Davide Furceri, Nikhil Patel, Shuvam Das

A novel quarterly database of discretionary fiscal measures is presented for 143 economies over 1950-2024. Policy passages are drawn from Economist Intelligence Unit country reports and processed in two steps. First, a rule-based parser identifies candidate text segments covering fiscal developments. Second, a large language model (GPT 4.1) classifies each episode by (i) net fiscal stance (expansion, contraction, or neutral) (ii) quantitative scale, and (iii) policy motivation. In the subsample that overlaps a prior study, agreement rates are 92 percent for the sign of policy actions and 90 percent for motivations behind the policies. The database extends narrative coverage by an order of magnitude in both the cross-section and the time dimension. Descriptive statistics document the distribution of tax and expenditure shocks across regions, income groups, and periods. As an illustration, country-specific Bayesian VARs yield four-quarter output multipliers for exogenous fiscal consolidations of 0.7 in the United States and 1.5 in Botswana, consistent with heterogeneous transmission predicted by theory.

C1010: **The signaling effects of fiscal announcements**

Presenter: **Francesco Zanetti**, University of Oxford, United Kingdom

Co-authors: Leonardo Melosi, Anna Rogantini Picco, Hiroshi Morita

Announcing a large fiscal stimulus may signal the government's pessimism about the severity of a recession to the private sector, impairing the stabilizing effects of the policy. Using a theoretical model, it is shown that these signaling effects occur when the stimulus exceeds expectations and are more noticeable during periods of high economic uncertainty. Analysis of a new dataset of daily stock prices and fiscal news in Japan supports these predictions. A method is introduced to identify fiscal news with different degrees of signaling effects and find that such effects weaken or, in extreme cases, even completely undermine the stabilizing impact of fiscal policy.

C1091: **Fiscal multipliers and economic risk**

Presenter: **Francesco Frangiamore**, University of Palermo, Italy

Co-authors: Davide Furceri, Domenico Giannone

The focus goes beyond average multipliers and state-dependent multipliers and extends the framework by analyzing the asymmetric effects of government expenditure through a location-scale (mean-variance) model. This approach allows studying the effects of government expenditure on the distribution of output without the need to pre-define the economic states. Results show that fiscal expansions have dual effects. They boost average economic activity (positive location effect) and reduce macroeconomic uncertainty (negative scale effect). This creates a compounding, stronger positive impact on downside risk. Fiscal policy does not just shift the output distribution rightward toward higher growth; it also compresses it with lower uncertainty, dramatically increasing the left tail where recession and negative growth scenarios occur. As a consequence, it is very powerful in periods when growth vulnerabilities are severe.

C1240: **Fiscal consolidations: Announcements and reality**

Presenter: **Roberto Perotti**, Bocconi University, Italy

Co-authors: Luca Sala

A large literature finds that consolidations based on announcements of tax increases lead to declines in GDP, while consolidations based on announcements of expenditure cuts are associated with almost no change in GDP. The purpose is to study the response of actual (as opposed to announced) discretionary government spending and revenues, and find that the former does not decline in what are classified (on the basis of announcements) as "expenditure-based" consolidations, but does decline by a large amount in "tax-based" consolidations; actual discretionary revenues move little in both regimes. Thus, opposite conclusions to those implied by most of the recent literature are reached.

CO292 Room BCB 208 ECONOMETRIC METHODS IN ENERGY, CLIMATE, AND RESOURCE RESEARCH

Chair: Christoph Wegener

C1194: **The impact of carbon pricing on electricity prices and volatility: evidence from the Italian power market**

Presenter: **Pierdomenico Dutillo**, University of Padova, Italy

Co-authors: Francesco Lisi

Carbon emissions are a major driver of global climate change, with the power sector among the largest contributors. Emissions from electricity generation depend on the balance between renewable and fossil fuel sources, which are influenced by both market dynamics and weather variability. Carbon pricing and emissions trading have emerged as central policy instruments to reduce emissions worldwide. Since its introduction in 2005, the European Union Emissions Trading Scheme (EU ETS) has required participating firms to surrender allowances for their greenhouse gas emissions,

with each allowance granting the right to emit one tonne of CO₂-equivalent. The aim is to examine the impact of carbon pricing on electricity prices and price volatility in the Italian power market. By modelling the conditional mean and variance of electricity price using parametric and non-parametric approaches, it demonstrates how carbon costs influence market behavior and shape electricity market dynamics.

C1160: **Disentangling oil market shocks using a proxy-FSVAR approach**

Presenter: **Linus Nuesing**, University of Konstanz, Germany

The purpose is to analyze the effect of the oil market on the U.S. economy through two different channels: the flow and the expectation channel. In particular, oil market shocks are divided into an oil flow supply and demand shock, as well as a supply and demand news shock. Furthermore, the interaction of these shocks with the macroeconomic and financial markets is examined. In order to identify the four oil price shocks, the Factor-SVAR approach is employed, and it is extended by a proxy equation that allows the identification of the two news shocks by external instruments. The factor structure in the residuals enables modelling of data containing different proxies for the same economic concept and the inclusion of variables that encompass measurement errors in their construction.

C1183: **Real options, commodities and natural resources under ambiguity**

Presenter: **Christian Ewald**, University of Glasgow, United Kingdom

Co-authors: Yihan Zou, Ankush Agarwal

The purpose is to explore real options in commodity and natural resource investments when decision-makers face not only risk but also ambiguity uncertainty about model parameters and future dynamics. Robust valuation and optimal timing frameworks are developed based on reflected backward stochastic differential equations (RBSDEs) combined with advanced Monte Carlo methods. In the context of large-scale commodity projects, parameter uncertainty in multi-factor models of spot prices, interest rates, and convenience yields generates a substantial ambiguity premium. Unlike risk, which tends to delay exercise, ambiguity often accelerates investment or harvesting decisions. This is illustrated through applications to aquaculture and forestry. For aquaculture projects, ambiguity in convenience yield estimation leads to earlier optimal harvesting and reduced project value. For forestry, both catastrophe risk, modeled as a Poisson process calibrated with wildfire data, and parameter uncertainty in lumber price dynamics are incorporated. This dual uncertainty reduces lease values and shifts optimal harvesting toward more conservative strategies, though carbon sequestration considerations can offset this effect by delaying harvests. Overall, the results highlight how ambiguity materially alters real options decisions in natural resource sectors and should be accounted for in project evaluation and climate-related risk management.

C1148: **Modelling the transition of fossil fuel corporations**

Presenter: **Isabel Figuerola-Ferretti**, Universidad Pontificia Comillas, Spain

Co-authors: Ioannis Paraskevopoulos

The aim is to develop a theoretical model to investigate how carbon pricing and the greenium influence the transition dynamics of energy companies. A stochastic model is developed, distinguishing economic stranding points, when fossil assets become unviable, and the full transition period, defined by the policy-determined time-frame for achieving climate goals. The model delivers explicit relationships between firm characteristics and regulatory conditions. Using quarterly panel data for 191 energy firms from 2015 to 2024, it is found that higher carbon prices and greenium levels significantly accelerate economic stranding of fossil-based activities. However, current economic signals from carbon pricing alone appear insufficient to meet 2050 net-zero targets. These findings highlight the need for targeted policies that enhance the transition of fossil fuel firms.

CO096 Room BCB 209 RECENT ADVANCES IN STATISTICAL FEW-SHOT LEARNING

Chair: Marta Nai Ruscone

C1396: **Analysis of Keystroke dynamics for multi-user systems**

Presenter: **Semhar Michael**, South Dakota State University, United States

Co-authors: Andrew Simpson

Keystroke dynamics has been used to both authenticate users of computer systems and detect unauthorized users who attempt to access the system. Monitoring keystroke dynamics adds another level to computer security, as passwords are often compromised. Keystrokes can also be continuously monitored long after a password has been entered and the user is accessing the system for added security. Many of the current methods proposed are supervised methods, assuming that the one authorized user of each keystroke is known a priori. This is not always true, for example, with businesses and government agencies, which have internal systems that multiple people have access to. This implies that unsupervised methods must be employed for these situations. A novel method is proposed that accounts for the lack of a one-to-one relationship between the number of users and the number of components, as well as accounts for known information based on when keystrokes were typed. Based on simulation studies and the motivating real-data example, the proposed model shows good performance.

C1410: **Induced likelihood-based methods for soft classification of u-processes arising in few shot learning problems**

Presenter: **Christopher Saunders**, South Dakota State University, United States

Co-authors: Janean Hanka, Danica Ommen, Semhar Michael

Few-shot learning problems are the most common class of pattern recognition tasks commonly encountered in forensic source identification. The data resources for this problem are commonly structured in a way that involves a large number (denoted as C) of classes (typically referred to as sources), each with two to five exemplars (typically referred to as control samples or knowns) per class of interest. Typically, a metric for measuring the dissimilarity between a pair of objects is constructed or learned from earlier studies. This learned metric is used to make predictions as a pseudo-metric for a k -nearest-neighbor classifier or as a metric for a kernel density estimator; neither of which leads to a probabilistic interpretation of the final output. Furthermore, the limited sample sizes within each class limit the ability to use a brute force method such as kernel density. To work around these structural constraints, a clustering-based method is used for the empirical distribution of pairwise scores within each class, which allows the pooling within distributions of scores across sources in a statistically rigorous manner to make stable probabilistic statements.

C1411: **Source inference for complex forensic evidence using a contrastive learning framework**

Presenter: **Danica Ommen**, Iowa State University, United States

Co-authors: Samuel Fox, Christopher Saunders, JoAnn Buscaglia

To interpret the value of forensic evidence resulting from paired item data, the common source identification framework asks: Do the items share a common unknown source, or do they come from two different unknown sources? This question can be addressed using a variety of forensic statistics techniques, including the usual two-stage, likelihood ratio, and Bayes factor approaches. Contrastive learning methods address the question using two major components: a method for quantifying the similarity (or dissimilarity) of pairs of evidence items, and a method for determining the best separation of within-source or between-source comparisons. Contrastive learning methods are particularly useful when the data derived from the evidence is high-dimensional or complex. In this case, the contrastive learning algorithms take advantage of high-performing artificial intelligence and machine learning tools to avoid specifying complicated probability models for the usual forensic statistics approaches. A contrastive learning algorithm framework is developed for complex evidence. The output of the contrastive learning algorithm can be used in a score-based likelihood ratio to interpret the value of evidence. Additional work is necessary to apply the method to the specific source question (whether an item came from a specific known source).

C1415: **On discriminant analysis for the statistical few-shot learning of skewed data**

Presenter: **Yana Melnykov**, University of Alabama, United States

Co-authors: Semhar Michael, Andrew Simpson

Few-shot learning with many classes and limited observations per class poses challenges, especially in high-dimensional settings. Traditional classification methods like linear or quadratic discriminant analysis assume that each class follows a multivariate normal distribution. However, quadratic discriminant analysis may suffer from unstable covariance estimates due to insufficient class-specific data. A recently developed approach based on clustering covariance matrices allows more flexible information sharing across classes. To relax the assumption of normality, a transformation-based procedure capable of handling skewed data is proposed. Monte Carlo simulations show promising results in parameter estimation as well as the classification accuracy.

CO301 Room BCB 211 NONCAUSAL ECONOMETRICS

Chair: Arthur Thomas

C0743: Disentangling drivers of EU allowance prices: A mixed causal and non-causal time series approach

Presenter: **Eduardo Serrubeco Marques**, University Paris Dauphine, France

Co-authors: Arthur Thomas, Olivier Massol

The purpose is to analyze the drivers of EU Allowance (EUA) price fluctuations within the EU emissions trading system (EU ETS), focusing on a methodology tailored to the carbon markets' unique characteristics. An extended autoregressive distributed lag (ARDL) model is employed, following prior studies, which allows for estimation of both short- and long-run effects irrespective of the integration order of the variables. The model captures the impact of key fundamental drivers' natural gas, coal, and Brent crude oil prices on EUA prices, accounting for nonlinearities using partial sum decompositions as proposed by another study. Fuel-switching dynamics between coal and gas are modeled based on deviations from the theoretical switch price. To address heteroscedasticity in the residuals, a mixed causal and non-causal framework is incorporated, and importance sampling-based partial filtering is applied to decompose EUA prices into fundamental (causal) and speculative (non-causal) components. The causal component is regressed on fundamentals using the NARDL model, successfully mitigating heteroscedasticity. The non-causal component is further analyzed via a VAR model to separate the effects of rational expectations from market sentiment, providing a comprehensive view of carbon price formation mechanisms.

C0449: Time aggregation of mixed causal noncausal autoregressive models with real and complex conjugate roots

Presenter: **Tomas del Barrio Castro**, University of the Balearic Islands, Spain

The aim is to study systematic and cumulant aggregation of mixed causal noncausal autoregressive models with alpha-stable innovations. It is shown that aggregation preserves noncausality; however, it may affect the modulus of the roots of both the causal and noncausal parts, as well as the frequency allocation of the roots (aliasing effect). In some cases, under cumulant sampling, pairs of complex conjugate factors or negative roots may not only shift from one frequency to another, but may even vanish entirely. The effects of aggregation on the fat-tail behavior of alpha-stable innovations are also analyzed after aggregation. The Monte Carlo simulations support the analytical findings

C0403: Prediction of bubbles in presence of alpha-stable aggregates moving averages

Presenter: **Arthur Thomas**, Paris Dauphine University - PSL, France

Co-authors: Gilles De Truchis, Sebastien Fries

Financial markets frequently exhibit boom-and-bust cycles that are incompatible with standard linear time series models. While anticipative heavy-tailed linear processes offer a promising alternative for modeling such phenomena, they impose uniform bubble patterns across different episodes, contradicting empirical evidence. A new model is introduced based on alpha-stable moving average aggregates that accommodates heterogeneous bubble dynamics. The theoretical properties of this model are established, demonstrating that it admits a semi-norm representation on a unit cylinder, thereby enabling the prediction of extreme trajectories with varying growth dynamics. A minimum distance estimation procedure is developed based on the joint characteristic function, and its asymptotic properties are established. Monte Carlo simulations confirm the estimator's good finite-sample performance across various specifications. The empirical application to the CBOE Crude Oil ETF Volatility Index successfully decomposes observed volatility dynamics into distinct components with different persistence properties, revealing that what appears as a single bubble episode actually consists of multiple superimposed processes with heterogeneous growth rates and crash probabilities.

C0404: Multivariate seminorm representation for alpha-stable moving average processes and path prediction

Presenter: **Gilles De Truchis**, University of Orlaans, France

Co-authors: Arthur Thomas, Sebastien Fries

For a univariate process $X(t)$ modeled as a two-sided alpha-stable moving average, the dependence between the past and future components of vectors of the form $(X(t-m), \dots, X(t), X(t+1), \dots, X(t+h))$, where $m \geq 0$ and $h \geq 1$, is encoded in their spectral measures. This dependence can be represented on unit cylinder sets for an appropriate seminorm only if the process is "anticipative enough." This framework allows for the explicit derivation of the conditional distribution of future paths when only the first $m+1$ components are observed and large in norm. From this, one can deduce a forecasting procedure capable of determining the crash date of an extreme event if the underlying univariate process is noncausal. This representation is extended to general multivariate alpha-stable moving averages, where new properties emerge in higher dimensions, where, in particular, the presence of a non-anticipative component does not rule out the existence of adequate seminorm representations. This leads to an interesting result: Extreme events (peak and crash dates) of purely causal processes (which are otherwise unpredictable) can be inferred from a noncausal variable estimated within the same multivariate system, as they are linked by the error term. This serves as proof of early warning. For practical applications, a closed-form expression is proposed for the conditional distribution of future paths in an alpha-stable mixed causal VAR model.

CO321 Room BCB 212 HIGH-DIMENSIONAL AND NONPARAMETRIC METHODS FOR PANEL DATA MODELS

Chair: Alexandra Soberon

C0558: Nonparametric estimation of smooth coefficients in fixed-effect panel data models

Presenter: **Taining Wang**, Capital University of Economics and Business, China

Co-authors: Feng Yao, Jun Cai

A kernel-based nonparametric estimator is proposed for a smooth coefficient panel data model with fixed effects. Without requiring a zero sum of fixed effects, an estimator is proposed that is easy to construct and computationally efficient. Eliminating the fixed effects through a local within transformation, a local linear estimation is performed for the coefficient functions associated with time-varying variables. The intercept coefficient function is further estimated, if present, through a difference of kernel weighted averages. The estimator's asymptotic properties are characterized under a large- n and large- T framework. It is demonstrated that the estimator is not asymptotically equivalent to the standard kernel estimator that ignores fixed effects. Through extensive simulation studies, the estimator's encouraging numerical performance and computational advantages are highlighted over existing kernel estimators in the literature. The empirical applicability is showcased by investigating a varying coefficient version of the environmental Kuznets curve through a panel of OECD countries.

C0725: Nonparametric identification in correlated random effects models

Presenter: **Daniel Henderson**, University of Alabama, United States

Co-authors: Daniel Henderson, Andros Kourtellis

A semiparametric estimation strategy is used to identify the unknown effects in a correlated random effects model. The estimation strategy allows

for more robust estimation in the presence of unobserved individual effects. Simulations and an empirical example show the methods work well in finite samples.

C1029: Weak instrumental variables due to nonlinearities in panel data: A super learner control function estimator

Presenter: **Monika Avila Marquez**, University of Geneva, Switzerland

A triangular structural panel data model is proposed with additive separable individual-specific effects to estimate the causal effect of a covariate on an outcome in the presence of unobservable confounders, some of which are time-invariant. In this context, a linear reduced-form equation may be problematic when the conditional mean of the endogenous covariate and instrumental variables is nonlinear, as ignoring such nonlinearity can lead to weak instruments. To address this, a triangular simultaneous equation model with a linear structural equation and a nonlinear reduced-form equation is introduced. The parameter of interest is the structural coefficient on the endogenous variable. Identification is achieved under standard exclusion restrictions using a control function approach. An estimator called the super learner control function estimator (SLCFE) is developed. The method involves two main steps and sample splitting across the individual dimension. First, the control function is estimated using a super learner; then, this estimate is used to correct for endogeneity in the structural equation. Consistency of the estimator is established, and its performance is evaluated through Monte Carlo simulations. Results show that SLCFE significantly outperforms conventional Within 2SLS estimators, especially when the reduced-form relationship is nonlinear.

C1043: Inference in high-dimensional two-way panel data models

Presenter: **Juan Manuel Rodriguez-Poo**, Universidad de Cantabria, Spain

Co-authors: Alexandra Soberon, Lindes Dominguez Diaz

The aim is to develop a consistent and asymptotically normal estimator for a triangular simultaneous two-way high-dimensional panel data model. The estimator addresses endogeneity arising from both individual and time effects, as well as the dependence between covariates and error terms. A two-stage procedure is proposed: First, removing fixed effects; second, applying instrumental variable estimation using regularization methods (lasso, cluster-lasso, and post-lasso) suitable for high-dimensional settings. Monte Carlo simulations demonstrate the estimator's favorable properties in terms of bias and RMSE, particularly when the regularization parameter is selected via cross-validation. Theoretical properties and practical implementation details are discussed, and the performance of the estimator is evaluated through extensive simulation studies.

CO144 Room BCB 213 ADVANCES IN PROBABILISTIC FORECASTING

Chair: Lukas Bauer

C0394: Shift-dispersion decompositions of Wasserstein and Cramer distances

Presenter: **Johannes Resin**, Goethe University Frankfurt, Germany

Co-authors: Timo Dimitriadis, Johannes Bracher, Daniel Wolffram

Divergence functions are measures of distance or dissimilarity between probability distributions that serve various purposes in statistics and applications. Decompositions of Wasserstein and Cramer distances are proposed, which compare two distributions by integrating over their differences in distribution or quantile functions, into directed shift and dispersion components. These components are obtained by dividing the differences between the quantile functions into contributions arising from shift and dispersion, respectively. The decompositions add information on how the distributions differ in a condensed form and consequently enhance the interpretability of the underlying divergences. It is shown that the decompositions satisfy a number of natural properties and are unique in doing so in location-scale families. The decompositions allow deriving sensitivities of the divergence measures to changes in location and dispersion, and they give rise to weak stochastic order relations that are linked to the usual stochastic and dispersive order. The theoretical developments are illustrated in two applications, where the focus is on forecast evaluation of temperature extremes and on the design of probabilistic surveys in economics.

C1231: Assessing the conditional calibration of interval forecasts using decompositions of the interval score

Presenter: **Sam Allen**, Karlsruhe Institute of Technology, Germany

Co-authors: Julia Burnello, Johanna Ziegel

Forecasts for uncertain future events should be probabilistic. Probabilistic forecasts are commonly issued as prediction intervals, which provide a measure of uncertainty in the unknown outcome whilst being easier to understand and communicate than full predictive distributions. The calibration of a prediction interval can be assessed by checking whether the probability that the outcome falls within the interval is equal to the nominal coverage level. However, such coverage checks are typically unconditional and therefore relatively weak. Although this is well known, there is a lack of methods to assess the conditional calibration of interval forecasts. The purpose is to demonstrate how this can be achieved via decompositions of the well-known interval (or Winkler) score. Notions of calibration are studied for interval forecasts, and a decomposition of the interval score is then introduced based on isotonic distributional regression. This decomposition exhibits many desirable properties, both in theory and in practice, which allow users to accurately assess the conditional calibration of interval forecasts. This is illustrated on simulated data and in three applications to benchmark regression datasets.

C0412: Kullback-Leibler-based characterizations of score-driven updates

Presenter: **Ramon de Punder**, University of Amsterdam, Netherlands

Score-driven models have been applied in some 400 published articles over the last decade. Much of this literature cites the optimality result in a prior study, which, roughly, states that succinctly small score-driven updates are unique in locally reducing the Kullback-Leibler divergence relative to the true density for every observation. This is at odds with other well-known optimality results; the Kalman filter, for example, is optimal in a mean-squared-error sense, but occasionally moves away from the true state. It is shown that score-driven updates are, similarly, not guaranteed to improve the localized Kullback-Leibler divergence at every observation. The seemingly stronger result in the prior study is due to their use of an improper (localized) scoring rule. Even as a guaranteed improvement for every observation is unattainable, it is proven that succinctly small score-driven updates are unique in reducing the Kullback-Leibler divergence relative to the true density in expectation. This positive, albeit weaker, result justifies the continued use of score-driven models and places their information-theoretic properties on a solid footing.

C0508: Online copula additive models for location, scale, and shape

Presenter: **Christian Schulz**, University of Duisburg-Essen, Germany

Co-authors: Simon Hirsch, Florian Ziel, Christoph Hanck

Large-scale streaming data are increasingly common in modern forecasting applications, particularly in the energy and finance sectors, and have motivated the development of online learning algorithms. In many empirical settings, jointly modeling two or more response variables conditional on covariates is of substantial interest. To achieve maximal flexibility, a generalized additive copula framework is adopted that models the marginal distributions and the copula separately, rather than assuming a multivariate distribution for the responses. Existing approaches are extended by incorporating an efficient online learning algorithm with exponential forgetting, based on online coordinate descent and LASSO-type regularization. The approach is validated in a forecasting study focused on the joint prediction of oil and gas prices. The proposed algorithms are implemented in the computationally efficient Python package `ondil`.

CO031 Room BCB M202 RECENT ADVANCES IN STATISTICAL FINANCE

Chair: Sayantan Banerjee

C0310: Sector-specific interrelationship between capital structure and sales growth: A Bayesian machine learning approach

Presenter: **Kousik Guhathakurta**, Indian Institute of Management Indore, India

Co-authors: Sujay Mukhoti

The aim is to investigate the sector-specific dynamics between capital structure and firm growth using a Bayesian machine learning framework. Traditional models often impose rigid parametric structures and struggle with endogeneity and model uncertainty. In contrast, the approach flexibly accommodates structural breaks, nonlinearities, and heterogeneous effects across industries by allowing the model space itself to evolve with the data. Findings reveal stark inter-sectoral differences in the leverage-growth nexus. In certain industries, higher financial leverage is positively associated with sustained sales growth, likely due to greater scale economies and access to external financing. Conversely, firms in some other sectors exhibit optimal growth under moderate leverage levels, beyond which excessive debt constrains operational flexibility and performance. These patterns highlight the importance of aligning financial strategies with sectoral characteristics such as capital intensity and competitive dynamics. Results provide a novel perspective for corporate managers, investors, and policymakers, emphasizing the need for industry-contingent capital structure decisions. The Bayesian machine learning framework offers a powerful and generalizable tool for understanding complex financial relationships in diverse economic contexts.

C0311: PDx: Adaptive credit risk forecasting model in digital lending using machine learning operations

Presenter: Sayantan Banerjee, Indian Institute of Management Indore, India

The aim is to present PDx, an adaptive, MLOps-driven system for forecasting credit risk using probability of default (PD) modeling in digital lending. Traditional PD models focus on accuracy at the development stage using complex ML algorithms, but often fail to adapt to evolving borrower behavior, leading to static models that degrade in production. Many lenders also struggle to deploy and maintain ML models effectively. PDx addresses these issues with a dynamic, end-to-end model lifecycle framework that includes continuous monitoring, automated retraining, and validation through a robust MLOps pipeline. A key innovation is a champion-challenger architecture, enabling regular updates and recalibration using recent data, ensuring resilience to data drift and shifting credit patterns. Empirical results show that ensemble tree models outperform others in default classification but require frequent updates to maintain performance. In contrast, logistic regression and neural networks exhibit faster degradation. By mitigating model decay and value erosion, PDx is especially effective for short-term, small-ticket digital loans where borrower behavior shifts quickly. PDx is validated using peer-to-peer, business, and auto loan datasets, demonstrating its scalability and adaptability for modern credit risk modeling.

C0715: A random projection based technique for change point detection in high-dimension

Presenter: Nilabja Guha, University of Manchester, United Kingdom

Co-authors: Jyotishka Datta

In many applications, such as economics, social science, and finance, changes in data-generating distributions are observed with time. The observed variable may depend on covariates through a mean structure, where the mean structure may change with time. There can also be changes in the underlying covariance structure. A Bayesian framework of change point estimation is presented for high-dimensional observations. Such high-dimensional observations may appear in many practical applications where the high-dimensional mean parameter or the covariance structure changes with time, such as high-frequency financial data. A lower-dimensional embedding is presented based on random projection. Change point estimation consistency and convergence rate are established even when the dimension of the observations can be much larger than the number of observations.

C0864: Bridging search behavior and market dynamics: A hybrid model for high-frequency financial data

Presenter: Soudeep Deb, Indian Institute of Management Bangalore, India

Co-authors: Archi Roy, Anitha Pathlavath

With the rapid evolution of financial markets, cryptocurrency has emerged as a unique and highly volatile asset class. The purpose is to examine how Google search volumes related to Bitcoin can serve as indicators of market sentiment and aid in forecasting price movements. Given the non-linearity and temporal dependencies in high-frequency financial data, a hybrid model combining non-parametric regression (NR) with long short-term memory (LSTM) networks is proposed. The NR-LSTM model outperforms traditional approaches, effectively capturing complex patterns and highlighting the value of search data in predicting Bitcoin price dynamics. These findings contribute to advancing cryptocurrency forecasting methodologies.

CO290 Room BCB 308 ADVANCES IN STATISTICAL GENOMICS

Chair: Jiebiao Wang

C0568: MODE: High resolution digital dissociation with deep multimodal autoencoder

Presenter: Qian Li, St. Jude Children's Research Hospital, United States

In single-cell biology, the complexity of tissues may hinder lineage cell mapping or tumor microenvironment decomposition, requiring digital dissociation of bulk tissues. Many deconvolution methods focus on transcriptomic assay, which is not easily applicable to other omics due to ambiguous cell markers and reference-to-target difference. The aim is to present MODE, a multimodal autoencoder pipeline linking multi-dimensional features to jointly predict personalized multi-omic profiles and cellular compositions, using pseudo-bulk data constructed by internal non-transcriptomic reference and external scRNA-seq data. The pseudo-bulk training data was generated by Dirichlet and Poisson distributions, using parameters initialized from the target data with a joint nonnegative matrix factorization (JNMF) model. MODE was evaluated through rigorous simulation experiments and real multi-omic data from multiple tissue types, outperforming nine deconvolution pipelines with superior generalizability and fidelity.

C0670: Supervised dimensionality reduction for visualization and prediction of microbiome data

Presenter: Rebecca Deek, University of Pittsburgh, United States

Advances in high-throughput sequencing technology have enabled researchers to directly measure microbial compositions. The resulting data are high-dimensional, with many microbes being observed in a single sample. Accordingly, dimensionality reduction algorithms, such as principal coordinates analysis, are commonly applied to microbiome data for visualization and feature extraction. Such algorithms are unsupervised and therefore can result in visualizations that fail to differentiate between distinct outcome groups and limit their utility beyond data visualization or exploration. As such, a supervised and covariate-adjusted principal coordinates analysis algorithm is proposed that incorporates similarities in the outcome, as well as the microbial compositions, into the reduced dimension data while simultaneously removing unwanted nuisance covariate effects. The method provides enhanced visualization and can be used for downstream modeling. The performance of the proposed method is illustrated using simulations and real microbiome data sets.

C1082: Heterogeneous causal mediation analysis using Bayesian additive regression trees

Presenter: Xu Qin, University of Pittsburgh, United States

Co-authors: Jiebiao Wang, Chen Liu

Causal mediation analysis provides insights into the mechanisms through which treatments affect outcomes. While mediation effects often vary across individuals, most existing methods focus solely on population-average effects, overlooking individual-level heterogeneity. To address this limitation, a Bayesian regression tree ensemble method is proposed that flexibly models non-linear relationships and captures treatment-by-mediator interactions in the mediation process. Using hierarchical posterior sampling, the approach provides credible intervals with nominal coverage rates for testing heterogeneous mediation effects. Additionally, regression tree summaries are leveraged to identify subgroups with distinct mediation effects, and SHapley Additive exPlanation (SHAP) values are employed to highlight key moderators and their influence on the mediation process.

Comprehensive simulations demonstrate the method's accuracy in estimating and inferring heterogeneous mediation effects. Finally, the method is applied to investigate the heterogeneous mediation of Alzheimer's disease pathology burden on the effect of apolipoprotein E (APOE) genotype on late-life cognition.

C1163: Recovering Isoform-level transcriptomics from sparse long-read single-cell RNA sequencing data

Presenter: **Wei Chen**, University of Pittsburgh, United States

Long read single-cell RNA sequencing (IscRNA-seq) quantifies single-cell isoform-level expression. Because gene-level signals split into many isoforms, feature numbers rise while counts fall, producing extreme sparsity. The resulting extremely sparse data represent highly partial observations of the transcriptome profiles, limiting the effectiveness of downstream statistical and bioinformatics methods. A graph-based diffusion method is presented to refine isoform expression per cell. Cell-cell similarity is computed using either gene-level or isoform-level expression data, by blending adapted SimRank similarity metrics with Gaussian-kernel weights from distances in a reduced space; a Markov process then propagates information across the similarity graph to estimate each cell's underlying transcriptomic profile. The framework preserves cell-type-specific signatures by limiting diffusion in the Markov process with adaptive neighborhood selection. Using simulated and real datasets, the method is shown to counteract sparsity, and biological information is recovered. Prior to refinement, important biological information, such as isoform correlations and differential expression along pseudotime, was largely hidden by dropout. After refinement, clear isoform correlations and changes along pseudotime emerge, improving interpretability. Overall, the approach enhances IscRNA-seq by leveraging cell graphs to recover biological information and reveal isoform-level expression patterns.

CO199 Room BCB 310 RECENT ADVANCES IN MODEL SELECTION

Chair: Marco Ferreira

C1177: Model selection for Bayesian dynamic clustering factor models

Presenter: **Tsering Dolkar**, Virginia Tech, United States

Co-authors: Marco Ferreira, Allison Tegge

The purpose is to consider model selection for Bayesian dynamic clustering factor models (BDCFM). BDCFM are novel models that combine latent factor models and hidden Markov models (HMMs) for the analysis of multivariate longitudinal data. As such, BDCFM perform concomitant dimension reduction, clustering, and estimation of the dynamic transitions of subjects among clusters. To select the number of clusters and the number of factors in BDCFM, an information criterion inspired by the Bayesian information criterion is proposed. To implement the information criterion, the computation of the likelihood for HMMs is extended using a forward-backward recursion to BDCFM. A simulation study shows that, when compared to competing approaches, the approach has favorable performance. The utility of the approach is shown with an application to a longitudinal study on recovery from substance use disorder.

C1146: Objective Bayes model selection for the Behrens-Fisher problem

Presenter: **Marco Ferreira**, Virginia Tech, United States

The aim is to propose a novel objective Bayes model selection approach for the Behrens-Fisher problem. The Behrens-Fisher problem concerns testing the equality of means of two populations when their variances are different. Typically, the data analyst would first test for the equality of the variances and, if the null hypothesis of such equality was rejected, the analyst would apply a test of equality of the means assuming unequal variances. Instead of this two-step approach, an omnibus objective Bayes solution that at the same time tests equality of means and of variances is developed. Simulation studies provide a comparison of the performance of our solution versus the performance of previously proposed methods. Applications to publicly available datasets illustrate the usefulness of the proposed approach.

C1150: What is in the model? A comparison of variable selection criteria and model search approaches

Presenter: **Shuangshuang Xu**, Virginia Tech, United States

Co-authors: Marco Ferreira, Allison Tegge

Variable selection methods are crucial for screening and identifying the most associated regression variables to a dependent variable of interest. In particular, a parsimonious model helps interpretation. Among the plethora of variable selection methods, BIC, AIC, and LASSO are three of the most widely used. The aim is to provide a comprehensive comparison among these three methods through a simulation study. For small numbers of regressors, the AIC and BIC are implemented with an exhaustive search of the model space. For a large number of regressors, the model space is explored using a genetic algorithm. The simulation study considers variable selection in linear models and generalized linear models. The results show that when compared with AIC and LASSO, for both small and large number of regressors, the BIC provides better performance in terms of correct identification rate and false discovery rate, while still having high recall.

C1449: Choice of number of factors and clusters in Bayesian clustering factor models

Presenter: **Allison Tegge**, Virginia Tech, United States

Co-authors: Marco Ferreira, Hwasoo Shin

A framework for concomitant dimension reduction and clustering based on a novel class of Bayesian clustering factor models (BCFM) was previously introduced. BCFM assume a factor model structure where the vectors of common factors follow a mixture of Gaussian distributions. The aim is to propose an information criterion to select the number of clusters and the number of factors. When compared to previously published competitor methods, this information criterion has favorable performance in terms of the correct selection of the number of clusters and the number of factors. Finally, the application of the information criterion is illustrated with a BCFM analysis of health care data.

CO224 Room BCB 311 INNOVATIONS IN STATISTICAL METHODS FOR COMPLEX DEPENDENCE

Chair: Honglang Wang

C0448: Rank-based inference for individualized treatment rules in single-index varying coefficient model

Presenter: **Honglang Wang**, Indiana University Indianapolis, United States

Co-authors: Yishan Cui

Individualized treatment rules (ITRs) provide critical guidance for patients by tailoring treatment decisions based on their specific covariates. However, deriving inferences for ITRs can be challenging, particularly when interactions between treatment and covariates are modeled non-parametrically, as this can introduce significant bias in estimating the ITRs. A unified rank-based inference procedure is proposed for ITRs under a semi-parametric single-index varying coefficient model, where the non-parametric coefficient function is assumed to be monotone increasing. To avoid direct estimation of the non-parametric function, the approach leverages maximum rank correlation. For hypothesis testing, the asymptotic distribution of the proposed estimator is not only derived using de-biasing techniques, but also the jackknife empirical likelihood is leveraged to test the significance of the treatment rule. The finite-sample performance of the proposed method is assessed through Monte Carlo simulations. The proposed method is further exemplified by its application to the ACTG175 data.

C0617: Composite likelihood inference for space-time point processes

Presenter: **Ganggang Xu**, University of Miami, United States

Co-authors: Abdollah Jalilian, Francisco Cuevas, Rasmus Waagepetersen

The dynamics of a rain forest are extremely complex, involving births, deaths, and growth of trees with complex interactions between trees, animals, climate, and environment. The patterns of recruits (new trees) and dead trees are considered between rainforest censuses. For the current census, regression models are specified for the conditional intensity of recruits and the conditional probabilities of death given the current trees

and spatial covariates. Regression parameters are estimated using conditional composite likelihood functions that only involve the conditional first-order properties of the data. When constructing assumption-lean estimators of covariance matrices of parameter estimates, mild assumptions of decaying conditional correlations in space are only needed, while assumptions regarding correlations over time are avoided by exploiting conditional centering of composite likelihood score functions. Time series of point patterns from rain forest censuses are quite short, while each point pattern covers a fairly big spatial region. To obtain asymptotic results, a central limit theorem is therefore used for the fixed timespan - increasing spatial domain asymptotic setting. This also allows handling the challenge of using stochastic covariates constructed from past point patterns. Conveniently, it suffices to impose weak dependence assumptions on the innovations of the space-time process.

C0712: Modularity-guided and dominant-set-based semi-supervised clustering with metric learning for functional data

Presenter: **Xiang Wang**, Shanxi University, China

Co-authors: Honglang Wang

Dominant-set clustering is a graph and game-theoretic approach that identifies cohesive groups by maximizing within-cluster similarity, different from common methods such as k-means, spectral, and hierarchical clustering. A dominant-set-based hierarchical bipartition procedure is proposed, formulated as a penalized optimization problem, with the tuning parameter selected to maximize the modularity of the resulting two clusters. The proposed method is applied to functional data clustering with a flexible choice of similarity measures between curves. It is not only robust to imbalanced groups but also to outliers, which overcomes the limitation of many existing clustering methods. A thorough semi-supervised clustering method is further proposed, which learns the metric by modularity maximization over a linear combination of similarity metric candidates from the labeled portion of the data, and performs hierarchical dominant-set based clustering tuned by modularity maximization. The proposed algorithm is not only able to learn a global metric but also able to learn individual metrics for each cluster, which permits innovative clustering with overlapping clusters. This is a general clustering method and is superiorly applicable to functional data, which in nature encompasses a variety of metrics for comparing curves. Empirical investigations using simulation studies and real data applications demonstrate the advantages of the proposed methods.

C1024: Spatial deconvolution and cell type-specific spatially variable gene detection in spatial transcriptomics

Presenter: **Yuehua Cui**, Michigan State University, United States

Co-authors: Haohao Su, Yuehua Cui

Spatial transcriptomics (ST) provides crucial insights into tissue-specific gene expression patterns in various cancer studies. Most ST data, such as that obtained from the 10x Visium platform, is captured at a spot resolution that measures gene expression across multiple cells, often originating from various cell types. Deconvolution of such multi-cellular data to infer cell type compositions is crucial for further downstream analysis. Recent methodological developments have greatly advanced the detection of spatially variable genes (SVGs), whose expression patterns are non-random across tissue locations. Given that many SVGs correlate with cell type compositions, a unified approach is introduced to identify both SVGs and cell type-specific SVGs (ctSVGs), integrated with ST deconvolution, under a linear mixed-effect model framework. The method, termed STANCE, ensures tissue rotation-invariant results, with a two-stage testing strategy: Initial SVG/ctSVG detection followed by ctSVG-specific testing. Its performance is demonstrated through extensive simulations and analyses of public datasets. Downstream analyses reveal STANCE's potential in spatial transcriptomics analysis.

CO262 Room BCB 312 STATISTICAL RESEARCH IN DESIGN OF EXPERIMENTS

Chair: Damianos Michaelides

C0397: Optimal robust designs with both centered and baseline factors

Presenter: **Xietao Zhou**, KCL, United Kingdom

Co-authors: Steven Gilmour

Traditional optimal designs are optimal under a pre-specified model. When the final fitted model differs from the pre-specified model, traditional optimal designs may cease to be optimal, and the corresponding parameter estimators may have larger variances. The Q_B criterion has been proposed to offer the capacity to consider hundreds of alternative models that could potentially be useful for data from a multifactor design. The Q_B criterion is extended to the scenario when eligible candidate models contain both baseline and centered parameterization factors. This shall be of interest in practice when some of the factors naturally do have a reasonable null state alongside other factors whose levels are equally important and are more naturally represented under the centered parameterization. The optimal designs are compared with their counterparts in the most recent literature and have shown that the projection capacity of eligible candidate the models/accuracy of estimation of models in terms of the A_γ criterion can be improved when the number of runs in the experiment is a multiple of 4 and have also examined and solved the same problem with no restrictions on the number of runs of the experiment so that it could be applied in a more general way in practice. The new version of the Q_B criterion, dealing with factors under both parameterizations, is presented, followed by evaluating the robust and accurate performance of the Q_B optimal designs found.

C0648: Searching for optimal design of experiments

Presenter: **Kalliopi Mylona**, King's College London, United Kingdom

Finding an optimal experimental design is computationally challenging, particularly in high-dimensional spaces or multi-objective problems. Real-world applications frequently involve competing design criteria that must be jointly balanced. Other optimal experimental designs may involve factors that have a functional nature. Examples can be found in the pharmaceutical industry, the environment, engineering, chemistry, and more. Algorithmic tools that provide a balance between computing speed and efficiency are presented, providing the framework for more reliable design processes.

C0243: Using Hasse diagrams to understand complex experimental designs

Presenter: **Simon Bate**, GlaxoSmithKline, United Kingdom

In many areas of scientific research, scientists routinely use complex experimental designs when conducting their experiments. With the advent of the modern statistical package, the scientist often carries out the analysis of data generated from such experiments. This can lead to incorrect and misleading results, especially if the scientist has failed to correctly identify the experimental design they are using and the effect the design has on the statistical analysis. The aim is to describe a procedure, and associated R package, that allows non-statisticians to identify the structure of the experimental design using a Hasse diagram. It is then possible to use this diagram, along with the randomization performed, to produce an appropriate model for the statistical analysis. Placing experimental design at the center of the statistical process not only reduces the complexity for non-statisticians but also improves the reliability of any statistical results generated.

C1336: Hasse diagrams and covariates

Presenter: **Hans-Michael Kaltenbach**, ETH Zurich, Switzerland

Co-authors: Damianos Michaelides, Simon Bate

Hasse diagrams provide a visual representation of the blocking and treatment structure of an experimental design and the treatment randomization. They allow automated specification of an adequate linear model and are a powerful nontechnical means for discussing design alternatives with non-specialists. In many experiments, additional information in the form of covariates is available and should be used in the analysis. However, covariate adjustments are usually seen as part of the analysis and not the design of an experiment, and hence have not yet been represented in Hasse diagrams. The aim is to present methods for integrating covariates into Hasse diagrams to represent the full ANOVA table, correctly account for covariates' degrees of freedom, and allow automated specification of analysis of covariance models. The use of covariate adjustments is showcased

in several standard designs of increasing complexity. It is also shown how this methodology can be extended to the modelling of continuous factors, thus allowing both categorical and continuous factors to be represented on Hasse diagrams.

CO078 Room BCB 313 EMERGING TOPICS IN STATISTICS AND DATA SCIENCE
Chair: Yichuan Zhao
C1044: Flexible hazards and cure models for dynamic prediction based on longitudinal biomarker measurements

Presenter: **Xuelin Huang**, University of Texas MD Anderson Cancer Center, United States

Co-authors: Can Xie, Ruosha Li, Nicholas Short, Christopher Flowers

To optimize personalized treatment strategies, treat patients efficiently, and extend their survival times, it is critical to accurately predict patients' prognoses at all stages, from disease diagnosis to follow-up visits. The longitudinal biomarker measurements during these visits are essential for this prediction purpose. Patients' ultimate concerns are cure and survival. However, in many situations, there is no clear biomarker indicator for cure. A comprehensive joint model of longitudinal and survival data is proposed, incorporating proportions of potentially cured patients. Formulas are provided for predicting an individual's probabilities of future cure and survival at any time point based on their current biomarker history. The survival distributions in the model are specified through flexible hazard functions with the proportional hazards as a special case, allowing other patterns such as crossing hazard and survival functions. Simulations show that, with these comprehensive and flexible properties, the proposed model outperforms commonly used models in terms of predictive performance, measured by the time-dependent area under the curve (AUC) of the receiver operating characteristic and the Brier score. The use and advantages of the proposed model are illustrated by its application to a study of patients with chronic myeloid leukemia.

C0454: A gradient boosting decision tree-based estimation method for the mixture cure model

Presenter: **Yingwei Peng**, Queen's University, Canada

Cure models are useful tools for analyzing censored survival data with a cured fraction. However, existing semiparametric estimation methods still rely on restrictive parametric assumptions, and existing nonparametric estimation methods only work with single covariates. A gradient boosting decision tree-based method is proposed to estimate the mixture cure model. The new method inherits the features of the original gradient boosting decision tree method, and more accurate estimates of the cure probability and the relative risk are provided for uncured subjects than existing methods when there are no a priori parametric assumptions on the forms of complex covariate effects in the model. This is demonstrated with small mean square errors in the estimates of the cure probability, relative risk score, and survival function in a simulation study with large samples. The method also has the potential to deal with high-dimensional covariates. The proposed method is illustrated with a large sample study of colon cancer.

C1306: On the comparison of paired proportions

Presenter: **Zhigang Zhang**, Memorial Sloan-Kettering Cancer Center, United States

When comparing proportions calculated from dependent binary data, one should take into account the correlation between the matched pairs. The purpose is to examine some commonly used statistics under new frameworks and assumptions. In addition, experimental designs and sample size estimation under such circumstances are discussed. Both simulation studies and real-world examples are presented for illustration.

C0557: Novel empirical likelihood method for the cumulative hazard ratio under stratified Cox models

Presenter: **Yichuan Zhao**, Georgia State University, United States

Evaluating the treatment's effect is a crucial topic in clinical studies. Nowadays, the ratio of cumulative hazards is often applied to accomplish this task, especially when those hazards may be nonproportional. The stratified Cox proportional hazards model, as an important extension of the classical Cox model, has the ability to flexibly handle nonproportional hazards. A novel empirical likelihood method is proposed to construct the confidence interval for the cumulative hazard ratio under the stratified Cox model. The large sample properties of the proposed profile empirical likelihood ratio statistic are investigated, and the finite sample properties of the empirical likelihood-based estimators under some different situations are explored in simulation studies. The proposed method was finally applied to perform statistical analysis on a real-world dataset on the survival experience of patients with heart failure.

CO243 Room BCB 402 RECENT METHODOLOGICAL ADVANCES IN BIOSTATISTICS
Chair: Zelalem Negeri
C0241: A general approach to modeling environmental mixtures with multivariate outcomes

Presenter: **Glen McGee**, University of Waterloo, Canada

Co-authors: Joseph Antonelli

An important goal of environmental health research is to assess the health risks posed by mixtures of multiple environmental exposures. In these analyses, flexible models like Bayesian kernel machine regression and multiple index models are appealing because they allow for arbitrary non-linear exposure-outcome relationships. However, this flexibility comes at the cost of low power, particularly when exposures are highly correlated and the health effects are weak, as is typical in environmental health studies. An adaptive index modeling strategy is proposed, that borrows strength across exposures and outcomes by exploiting similar mixture component weights and exposure-response relationships. In the special case of distributed lag models, in which exposures are measured repeatedly over time, co-clustering of lag profiles and exposure-response curves is jointly encouraged to more efficiently identify critical windows of vulnerability and characterize important exposure effects. The proposed approach is then extended to the multivariate index model setting where the true index structure is unknown, and variable importance measures are introduced to quantify component contributions to mixture effects. Using time series data from the National Morbidity, Mortality and Air Pollution Study, the proposed methods are demonstrated by jointly modeling three mortality outcomes and two cumulative air pollution measurements with a maximum lag of 14 days.

C0603: Some new advances in similarity-based patient-outcome predictive modeling

Presenter: **Joel Dubin**, University of Waterloo, Canada

Earlier work has shown that similarity-based predictive models can improve upon predictive performance, as compared to using the entire training data to help build models, particularly regarding model discrimination for binary responses. The focus is on the similarity-based modeling for joint consideration of model calibration and discrimination, as well as for dynamic prediction models. Properties of the methods are investigated in comprehensive simulation studies, and the methods are demonstrated through a separate analysis of a publicly available intensive care unit (ICU) database.

C0704: Estimating average treatment effects when treatment data are absent in a target study

Presenter: **Lan Wen**, University of Waterloo, Canada

Researchers are often interested in understanding the causal effect of treatment interventions. However, in some cases, the treatment of interest readily available in a randomized controlled trial is either not directly measured or entirely unavailable in observational datasets. This challenge has motivated the development of stochastic incremental propensity score interventions, which operate on post-intervention exposures affected by the treatment of interest, with the aim of approximating the causal effects of the treatment intervention. Yet, a key challenge lies in the fact that the precise distributional shift of these post-intervention exposures induced by the treatment is typically unknown, making it uncertain whether the approximation truly reflects the causal effect of interest. The primary objective is to explore data integration methodologies to characterize a distribution of post-intervention exposures resulting from the treatment in an external dataset, and to use this information to estimate counterfactual

mean outcomes under treatment interventions, in settings where the observational data lack treatment information and the external data may not contain measurements of the outcome of interest. The underlying assumptions required for this approach are discussed, and methodological guidance on estimation strategies is provided to address these challenges.

C0834: **Weighted vs unweighted multivariate bilinear regression models**

Presenter: **Jemila Hamid**, University of Ottawa, Canada

Inference for the growth curve models (GCMs) and the generalized multivariate analysis of variance (GMANOVA) models in general leads to bilinear projections. For this reason, the GCMs are referred to as multivariate bilinear regression models. Assuming multivariate normality, explicit likelihood solutions for the parameters of the GCMs exist. Both the least squares and the likelihood procedures lead to weighted inference, where the weight is the pooled sample variance-covariance matrix. Nevertheless, the variance-covariance estimator for such bilinear regression models is different from the pooled sample variance-covariance matrix; hence, other weighting considerations can be made. It is hypothesized that weights that take the bilinear nature of the design into consideration may, in fact, lead to better inference, at least in some situations. The aim is to present results from extensive simulations performed to examine the relative efficiency gained by using weighted bilinear regression compared to unweighted inference. The relative bias and relative efficiency of the available weighting strategies are also discussed, and a novel weighting algorithm is introduced, involving iteratively re-weighted likelihood estimators, which also accounts for the bilinear nature of the GCMs. Using simulations, it is demonstrated that this iteratively re-weighted approach leads to robust estimators in the presence of outliers and when the model is misspecified.

CO092 Room BCB 403 STATISTICAL ANALYSES OF COMPLEX AND MULTIPLEX/MULTILAYER NETWORKS

Chair: Fangzheng Xie

C0406: **Graph encoder embedding**

Presenter: **Cencheng Shen**, University of Delaware, United States

The purpose is to discuss the graph encoder embedding, its theoretical properties under random graph models, its scalability and numerical performance in vertex classification and clustering, as well as recent method advancements, including refined embedding and principal embedding.

C0408: **Random-walk debiased inference for contextual ranking model with application in large language model evaluation**

Presenter: **Yichi Zhang**, Statistics Department, Indiana University Bloomington, United States

A debiased inference framework is proposed to infer the ranking structure in the contextual Bradley-Terry-Luce (BTL) model. For the pairwise item comparison, a novel random-walk debiased estimator is introduced to efficiently aggregate the estimating functions of different item pairs. To approximate the nuisance preference score functions in the debiased estimator, a nonparametric maximum likelihood-based method is further introduced that can leverage many loss-minimization methods, e.g., the spline regression and deep neural networks. With decently estimated nuisance functions, the debiased estimator yields a tractable distribution and achieves the semiparametric efficiency lower bound asymptotically. The method is further extended for multiple hypothesis testing by incorporating the multiplier bootstrap techniques and adapting it to accommodate the distributional shift of contextual variables. Thorough numerical studies are provided to validate the statistical properties of the method, and its applicability is showcased in evaluating large language models based on human preferences under different contexts.

C0486: **Minority representation in network rankings: Methods for estimation, testing, and fairness**

Presenter: **Peter MacDonald**, University of Waterloo, Canada

Co-authors: Eric Kolaczyk, Hui Shen

Networks, composed of nodes and their connections, are widely used to model complex relationships across various fields. Centrality metrics often inform decisions such as identifying key nodes or prioritizing resources. However, networks frequently suffer from missing or incorrect edges, which can systematically affect centrality-based decisions and distort the representation of certain protected groups. To address this issue, a formal definition of minority representation is introduced, measured as the proportion of minority nodes among the top-ranked nodes. Systematic bias is modeled against minority groups by using group-dependent missing edge errors. Methods are proposed to estimate and detect systematic bias. Asymptotic limits of minority representation statistics are derived under canonical network models and used to correct the representation of minority groups in node rankings. Simulation results demonstrate the effectiveness of the estimation, testing, and ranking correction procedures, and the methods are applied to a contact network, showcasing their practical applicability.

C0544: **Efficient analysis of latent spaces in heterogeneous networks**

Presenter: **Yinqiu He**, University of Wisconsin - Madison, United States

A unified framework is proposed for efficient estimation under latent space modeling of heterogeneous networks. A class of latent space models is considered that decomposes latent vectors into shared and network-specific components across networks. A novel procedure is developed that first identifies the shared latent vectors and further refines estimates through efficient score equations to achieve statistical efficiency. Oracle error rates for estimating the shared and heterogeneous latent vectors are established simultaneously. The analysis framework offers remarkable flexibility, accommodating various types of edge weights under general distributions.

CO308 Room BCB 406 STATISTICAL METHODS IN NETWORKS AND HIGH-DIMENSIONAL DATA

Chair: Vince Lyzinski

C1330: **Integrate co-expression networks into multivariate regression for multi-omics analysis**

Presenter: **Shuo Chen**, University of Maryland, United States

Co-authors: Hwiyoung Lee

Accounting for dependence among high-dimensional variables in omics data analysis is critical to obtain accurate and reliable statistical inference. Although latent, omics variables often exhibit structured correlation/co-expression patterns. However, there are few methods explicitly accounting for such structured dependence in the statistical analysis of omics data. To address this methodological gap, the aim is to propose CoReg, which integrates co-expression patterns into multivariate regression analysis. Computationally efficient algorithms are developed to implement CoReg, and they are applied to extensive simulation studies and real-world omics data analyses. It is shown in simulations that CoReg substantially improves the accuracy of statistical inference and replicability across studies. These findings suggest that CoReg is well-suited for omics data analysis with dependence adjustment, analogous to how mixed-effects models handle repeated measures in lower-dimensional settings.

C1363: **CITE-ME: Controlling for induced triangles in estimating network model evolution**

Presenter: **Benjamin Leinwand**, Stevens Institute of Technology, United States

Co-authors: Keith Levin

In the standard degree corrected block model (DCBM), once the parameters are fixed, edges are conditionally independent of one another. This ignores effects like triadic closure observed in real networks. Conversely, many network models that include conditional edge dependence do not explicitly incorporate both degree correction features for each node and community structure. The purpose is to discuss a network generative model whereby edges are sampled one at a time without replacement. The initial sampling weights of the edges depend on the underlying affinities between the incident nodes absent any triadic closure, which may follow the structure of the DCBM probability matrix. As the sampling continues, though, edge sampling weights are updated to promote triadic closure within the network. The details of this generative mechanism and its effects on the resulting network, including distorting community structures, are discussed.

C1294: Solution diversification in graph matching matched filters*Presenter:* **Zhirui Li**, University of Pennsylvania, United States*Co-authors:* Ben Johnson, Daniel Sussman, Carey Priebe, Vince Lyzinski

The purpose is to present a novel approach for finding multiple noisily embedded template graphs in a very large background graph. The method builds upon the graph-matching-matched-filter technique proposed in a prior study, with the discovery of multiple diverse matchings being achieved by iteratively penalizing a suitable node-pair similarity matrix in the matched filter algorithm. In addition, algorithmic speed-ups are proposed that greatly enhance the scalability of the matched-filter approach. Theoretical justification of the methodology is presented in the setting of correlated Erdos-Renyi graphs, showing its ability to sequentially discover multiple templates under mild model conditions. The method's utility is additionally demonstrated via extensive experiments both using simulated models and real-world datasets, including human brain connectomes and a large transactional knowledge base.

C1407: Euclidean mirrors and first-order changepoints in network time series*Presenter:* **Avanti Athreya**, Johns Hopkins University, United States*Co-authors:* Tianyi Chen, Zachary Lubbets, Youngser Park, Carey Priebe

The purpose is to describe a model for a network time series whose evolution is governed by an underlying stochastic process, known as the latent position process, in which network evolution can be represented in Euclidean space by a curve, called the Euclidean mirror. The notion of a first-order changepoint is defined for a time series of networks, and a family of latent position process networks is constructed with first-order changepoints. It is proven that a spectral estimate of the associated Euclidean mirror localizes these changepoints, even when the graph distribution evolves continuously, but at a rate that changes. Simulated and real data examples on brain organoid networks show that this localization captures empirically significant shifts in network evolution.

CO312 Room BCB 407 STATISTICAL INFERENCE AND LEARNING IN COMPLEX SYSTEMS**Chair: Doudou Zhou****C0611: Optimal assortment inference within an online learning framework***Presenter:* **Shuting Shen**, National University of Singapore, Singapore

The modern retailing system is witnessing fast updating in customer behaviors, entailing adaptive policies to capture the dynamics of customer preferences. To manage the risks associated with changing customer preferences, it is important to develop an online framework that quantifies the uncertainty of the optimal assortment adaptively. The combinatorial inference of the optimal assortment is studied under the contextual multinomial logit model. Customer choice outcomes are actively collected over a series of time points, where the contextual information includes embedding vectors that capture the customer-product dynamics and revenue parameters. The offer set is adaptively selected at each time based on historical data. An inferential procedure is proposed that constructs a discrete confidence set for the true optimal assortment, facilitating inference on key properties of the optimal assortment. The temporal dependency and combinatorial structure of the Hessian matrix create challenges for convergence analysis. To address these, new anti-concentration bounds are developed for the Gaussian maxima difference. Furthermore, the high dimensionality is addressed by employing discretization and subspace projection techniques. Theoretical guarantees are provided on both the validity and power of the inferential procedure, and information-theoretic lower bounds are established for the required signal strength, which match the upper bounds of the procedure up to logarithmic factors.

C0618: Root cause discovery via permutations and Cholesky decomposition*Presenter:* **Jinzhou Li**, Stanford University, United States*Co-authors:* Benjamin Chu, Julien Gagneur, Marloes Maathuis

The motivation is from the following problem: Can the disease-causing gene in a patient affected by a monogenic disorder be identified? This problem is an instance of root cause discovery. In particular, the aim is to identify the intervened variable in one interventional sample using a set of observational samples as reference. A linear structural equation model is considered where the causal ordering is unknown. It begins by examining a simple method that uses squared z-scores and characterizing the conditions under which this method succeeds and fails, showing that it generally cannot identify the root cause. It is then proven, without additional assumptions, that the root cause is identifiable even if the causal ordering is not. Two key ingredients of this identifiability result are the use of permutations and the Cholesky decomposition, which allows exploiting an invariant property across different permutations to discover the root cause. Furthermore, permutations that yield the correct root cause are characterized and, based on this, a valid method for root cause discovery is proposed. This approach is also adapted to high-dimensional settings. Finally, the performance of the methods is evaluated through simulations, and the high-dimensional method is applied to discover disease-causing genes in the gene expression dataset that motivates this work.

C1034: MATES: Multi-view aggregated two-sample test*Presenter:* **Zexi Cai**, Columbia University, United States*Co-authors:* Wenbo Fei, Doudou Zhou

The two-sample test is a fundamental problem in statistics with a wide range of applications. In the realm of high-dimensional data, nonparametric methods have gained prominence due to their flexibility and minimal distributional assumptions. However, many existing methods tend to be more effective when the two distributions differ primarily in their first and/or second moments. In many real-world scenarios, distributional differences may arise in higher-order moments, rendering traditional methods less powerful. To address this limitation, a novel framework is proposed to aggregate information from multiple moments to build a test statistic. Each moment is regarded as one view of the data and contributes to the detection of some specific type of discrepancy, thus allowing the test statistic to capture more complex distributional differences. The novel multi-view aggregated two-sample test (MATES) leverages a graph-based approach, where the test statistic is constructed from the weighted similarity graphs of the pooled sample. Under mild conditions on the multi-view weighted similarity graphs, theoretical properties of MATES are established, including a distribution-free limiting distribution under the null hypothesis, which enables straightforward type-I error control. Extensive simulation studies and real data analysis demonstrate that MATES effectively distinguishes subtle differences between distributions.

C1286: Towards the most powerful test statistics for randomization tests*Presenter:* **Zijun Gao**, University of Southern California, United States

Randomization tests are well celebrated for their model-lean validity, but their power is highly dependent on the choice of test statistic. The aim is to first characterize the most powerful test statistic, assuming oracle knowledge of the nuisance functions that characterize the data-generating process. A method is then proposed to estimate these nuisance functions without sample splitting. The resulting randomization test is valid and particularly powerful for small samples and imbalanced designs.

CO045 Room BCB 408 ADVANCES IN MONTE CARLO METHODS FOR BAYESIAN STATISTICS**Chair: Alexandros Beskos****C0269: An unadjusted Barker algorithm for high-dimensional sampling without M-smoothness***Presenter:* **Samuel Livingstone**, University College London, United Kingdom

The aim is to introduce a recently proposed skew-symmetric numerical scheme for stochastic differential equations. At each step, an innovation is generated by skewing a Gaussian random variable in the direction of the drift. If the level of skew is chosen appropriately it can be shown that the scheme accurately approximates the underlying diffusion process in both the weak and strong sense. Some results are shown about both the finite time and long time simulation, in the latter case applying the scheme to the overdamped Langevin diffusion. The scheme resembles an

unadjusted version of the Barker proposal Metropolis-Hastings algorithm. A key feature is that no Lipschitz requirement on the drift is needed in order to establish strong convergence in the mean-squared sense or for a well-behaved sampling algorithm that converges geometrically quickly to an equilibrium with controllable discretization error.

C1129: Multiscale perspectives on computational statistical methods

Presenter: Deniz Akyildiz, Imperial College London, United Kingdom

Multiscale perspectives are presented on a number of long-standing problems in computational statistics. First, we show how multiscale and stochastic averaging techniques can be utilized to understand certain practical implementations of statistical inference methods, such as expectation-maximization. These insights are then used to understand a popular training method for generative models, namely contrastive divergence. Finally, if time permits, a novel application of multiscale ideas to multimodal sampling problems is introduced. A high-level introduction is provided to multiscale processes based on stochastic differential equations and how these can be connected to a number of computational statistical ideas and beyond.

C1222: Diff-Fusion: Bayesian fusion via denoising diffusions

Presenter: Adrien Corenflos, University of Warwick, United Kingdom

Co-authors: Adam Johansen, Gareth Roberts

The Fusion problem is concerned with the distributed aggregation of independent subposteriors $(f^c(x_c))_{c=1}^C$ into a common target distribution $\pi(x) \propto \prod_{c=1}^C f^c(x)$. Several methods have been proposed to treat this problem, which can be broadly classified into two categories: Approximate solutions, which exhibit an irreducible bias, and exact solutions, which are usually computationally expensive. The aim is to propose a novel sampling algorithm for the Fusion, based on a denoising diffusion model, where the Fusion distribution is noised towards independence by a set of Langevin dynamics, and then denoised back to the target distribution via sequential Monte Carlo methods. The method is embarrassingly parallel and provides an asymptotically consistent approximation of the Fusion distribution at a fraction of the cost of prior alternatives. A theoretical analysis of the method is provided, showing that, when a thin discretization is used, the algorithm is close to the exact solution of the Fusion problem, and its performance is illustrated on a set of numerical experiment benchmarks.

C1502: Manifold Markov chain Monte Carlo methods for Bayesian inference in diffusion models

Presenter: Alexandros Beskos, University College London, United Kingdom

Co-authors: Alexandre Thiery, Matt Graham

Bayesian inference for nonlinear diffusions, observed at discrete times, is a challenging task that has prompted the development of a number of algorithms, mainly within the computational statistics community. We propose a new direction, and accompanying methodology - borrowing ideas from statistical physics and computational chemistry - for inferring the posterior distribution of latent diffusion paths and model parameters, given observations of the process. Joint configurations of the underlying process noise and of parameters, mapping onto diffusion paths consistent with observations, form an implicitly defined manifold. Then, by making use of a constrained Hamiltonian Monte Carlo algorithm on the embedded manifold, we are able to perform computationally efficient inference for a class of discretely observed diffusion models. Critically, in contrast with other approaches proposed in the literature, our methodology is highly automated, requiring minimal user intervention and applicable across a range of settings, including elliptic or hypo-elliptic systems; observations with or without noise; and linear or non-linear observation operators. Exploiting Markovianity, we propose a variant of the method with complexity that scales linearly in the resolution of path discretisation and the number of observation times.

CO187 Room BCB 409 RECENT ADVANCES IN CAUSAL ANALYSIS AND SPATIAL STATISTICS

Chair: Jonathan Bradley

C1328: Estimation of cluster-specific causal effects on spatially associated survival data using SoftBART

Presenter: Indrabati Bhattacharya, Florida State University, United States

Co-authors: Durbadal Ghosh, George Rust, Debajyoti Sinha

The aim is to propose a novel Bayesian approach to estimate causal effects in spatially clustered survival data. Using soft Bayesian additive regression trees (SBART), a nonparametric regression is introduced for a log-Normal survival model that accommodates spatial associations among unknown cluster effects through a directed acyclic graph autoregressive (DAGAR) model. A two-stage approach is employed, which entails estimating the propensity score in the first step and incorporating it as a confounder of the outcome model in the second step. In the simulation study, the method is compared with existing approaches under various simulation scenarios, including both correctly specified and misspecified outcome models, to demonstrate the superior performance of the method. The method is then applied to analyze the causal effect of treatment delay (TD) on post-treatment survival of breast cancer patients from the Florida Cancer registry (FCR). The analysis produces the county-specific as well as state-wide assessment of the causal effects while accommodating spatial association among counties.

C1360: A subjective Bayesian approach to address unchecked assumptions in a causal analysis

Presenter: Jaiyool Kim, Florida State University, United States

Co-authors: Jonathan Bradley, Indrabati Bhattacharya

Causal analyses typically rely on four standard assumptions, yet in many applications these assumptions may not hold. As a result, estimates may be biased and the uncertainty understated. The purpose is to introduce a subjective Bayesian framework that introduces a dichotomous parameter, which equals one when causal assumptions hold and equals zero otherwise. When the assumptions do not hold, the potential outcome is assumed to be unobserved and follows a semi-parametric statistical model, and the prior probability that the assumptions hold is pre-specified to a particular value p . This new perspective leads to a new inferential question. Namely, what is the smallest value of p that leads to a significant average treatment effect? In simulation studies, data is generated from an observational model that violates consistency through a spatially structured mediator, but intentionally fit a potential-outcome model that omits this mediator, mimicking an analysis with an unverifiable consistency assumption. Bias and interval coverage of the average treatment effect are evaluated. The method is also applied to county-level data on fine particulate matter exposure and COVID-19 mortality in the United States. The results demonstrate how to address unchecked subjective variability present due to unverifiable assumptions.

C1400: A method to incorporate subsampling into Bayesian models for high-dimensional spatial data

Presenter: Jonathan Bradley, Florida State University, United States

Spatial statistical models with weakly stationary process assumptions have become standard in spatial statistics. However, one disadvantage of such models is the computation time, which rapidly increases with the number of data points. The goal is to apply an existing subsampling strategy to standard spatial additive models and to derive the spatial statistical properties. The approach has the advantage that one does not require any additional restrictive model assumptions. That is, computational gains increase as model assumptions are removed when using the model framework. This provides one solution to the computational bottlenecks that occur when applying methods such as Kriging to big data. Several properties of this new approach are provided in terms of moments, sill, nugget, and range under several sampling designs. An advantage of the approach is that it subsamples without throwing away data, and can be implemented using datasets of any size that can be stored. The results of the spatial data subset model approach are presented on simulated datasets and on a large dataset consisting of 150,000 observations of daytime land surface temperatures measured by the MODIS instrument onboard the Terra satellite.

C1409: Stochastic gradient MCMC for massive geostatistical data*Presenter:* **Reetam Majumder**, University of Arkansas, United States*Co-authors:* Mohamed Abba, Brian Reich, Brandon Feng

Gaussian processes (GPs) are commonly used for prediction and inference for spatial data analyses. However, since estimation and prediction tasks have cubic time and quadratic memory complexity in the number of locations, GPs are difficult to scale to large spatial datasets. The Vecchia approximation induces sparsity in the dependence structure and is one of several methods proposed to scale GP inference. The purpose is to add to the substantial research in this area by developing a stochastic gradient Markov chain Monte Carlo (SGMCMC) framework for efficient computation in GPs. At each step, the algorithm subsamples a minibatch of locations and subsequently updates process parameters through a Vecchia-approximated GP likelihood. Since the Vecchia-approximated GP has a time complexity that is linear in the number of locations, this results in scalable estimation in GPs. Through simulation studies, SGMCMC is demonstrated to be competitive with state-of-the-art scalable GP algorithms in terms of computational time and parameter estimation. An application of the method is also provided using the Argo dataset of ocean temperature measurements.

CC399 Room BCB 210 ECONOMETRIC AND FINANCIAL MODELLING**Chair: Angeles Carnero****C0187: Towards modeling lifetime default risk: Exploring different subtypes of recurrent event Cox-regression models***Presenter:* **Tanja Verster**, North-West University, South Africa*Co-authors:* Arno Botha, Bernard Scheepers

In the pursuit of modeling a loan's probability of default (PD) over its lifetime, repeat default events are often ignored when using Cox proportional hazard (PH) models. Excluding such events may produce biased and inaccurate PD-estimates, which can compromise financial buffers against future losses. Accordingly, a few subtypes of Cox models that can incorporate recurrent default events are investigated. Using South African mortgage data, both the Andersen-Gill (AG) and the Prentice-Williams-Peterson (PWP) spell-time models are explored. These models are compared against a baseline that deliberately ignores recurrent events, called the time to first default (TFD) model. Models are evaluated using Harrell's c-statistic, adjusted Cox-Sell residuals, and a novel extension of time-dependent receiver operating characteristic (ROC) analysis. From these Cox models, it is demonstrated how to derive a portfolio-level term structure of default risk, which is a series of marginal PD estimates at each point of the average loan's lifetime. While the TFD and PWP models do not differ significantly across all diagnostics, the AG model underperformed expectations. Depending on the prevalence of recurrent defaults, one may therefore safely ignore them when estimating lifetime default risk. Accordingly, the current practice of using Cox-modeling is enhanced in producing timeous and accurate PD-estimates under IFRS 9.

C1106: Longevity risk hedging via non-Gaussian state-space mortality models: A mean-variance-skewness-kurtosis approach*Presenter:* **Wai-Sum Chan**, The Hang Seng University of Hong Kong, Hong Kong*Co-authors:* Johnny S-H Li

Longevity risk has recently become a high-profile risk among insurers and pension plan sponsors. One way to mitigate longevity risk is to build a hedge using derivatives that are linked to mortality indexes. Longevity hedging methods are often based on the normality assumption, considering only the variance but no other (higher) moments. However, strong empirical evidence suggests that mortality improvement rates are driven by asymmetric and fat-tailed distributions, so that existing longevity hedging methods should be expanded to incorporate higher moments. The purpose is to fill the gap by adopting a mean-variance-skewness-kurtosis approach based on non-Gaussian extensions of commonly used stochastic mortality models, formulated in a state-space setting. On the basis of a general representation of these models, the authors derive (approximate) analytical expressions for the moments of the present values of the hedging instruments and the liability being hedged. These expressions are then integrated with a polynomial goal programming model, from which the optimal hedge portfolio is identified. Finally, theoretical results are demonstrated with a real mortality data set and a range of hedger preferences.

C1285: Understanding asset class interdependence: A vector copula approach*Presenter:* **Ivan Medovikov**, Brock University, Canada

The purpose is to assess the structure of dependence between major asset classes, including equities, fixed income, real estate, and cryptocurrencies, and trace the evolution of these links through recent major market events such as COVID and Liberation Day tariffs. A measure of vector dependence of a prior study is used to separate dependence that exists within each asset class from dependence that exists between these asset classes as a whole, where each class is treated as a single vector of asset returns. It is therefore demonstrated how the concept of vector, or "group" dependence can provide a deeper understanding of higher-dimensional links that exist in the financial markets.

C1236: A simulation approach to robust risk management of derivative products*Presenter:* **Bertrand Tavin**, EMLYON Business School, France

The purpose is to consider the problem of assessing and hedging the risk carried by a portfolio of non-standard derivative products managed with a family of parametric models. The problem is first formalized, and the framework is defined in which it can be solved by using a constrained simulation approach with respect to model parameters. The approach is suitable for an agent who may be agnostic with respect to a prior model and who may wish to account for expert views on the range of possible model parameters. Instead of breaking them into several sub-problems, the proposed methodology has the important advantage of answering the risk measurement and robust hedging questions in one step. Namely, the agent needs to run the simulation just once to get the desired answers. Numerical results obtained with recent market data when applying the method to a portfolio of variance swaps and forward-start options valued with models incorporating stochastic volatility and jumps are presented. In addition, the method can easily accommodate additional features such as cardinality constraints for the robust hedging strategy or transaction costs.

CC438 Room BCB M201 FINANCIAL ECONOMETRICS II**Chair: Massimiliano Caporin****C1117: Cryptocurrency in war: A double-edged sword?***Presenter:* **Jeffrey Chu**, Renmin University of China, China*Co-authors:* Stephen Chan, Yuanyuan Zhang

The purpose is to examine the short-term impact of the Russia-Ukraine war on the high-frequency digital asset markets. An event study approach is applied, focusing on the initial months of the war, and hourly returns of cryptocurrencies, DeFi tokens, and metaverse tokens are analyzed. It is found that negative war-related events have both an immediate and sustained impact on cryptocurrencies and DeFi tokens, likely due to a series of negative events leading to positive returns. In contrast to stocks and commodities like gold, cryptocurrencies and DeFi tokens exhibit positive and significant cumulative returns following negative war-related events. This suggests that these assets could serve as diversifiers or hedges against such events, similar to the political property observed for oil. Importantly, these findings provide preliminary insights into the ongoing Russia-Ukraine conflict and help to understand the impact of military conflict on cryptocurrency markets more broadly.

C1176: The interplay between cryptocurrencies and phishing crimes*Presenter:* **Yuanyuan Zhang**, American University of Sharjah, United Arab Emirates*Co-authors:* Stephen Chan, Jeffrey Chu

The purpose is to explore the relationship between global phishing crimes and cryptocurrency usage, with a specific focus on key metrics of

the Ethereum platform. Utilizing sample monthly data covering the period of January 2016 to December 2022, inclusive, a quantile on quantile regression approach is employed to analyze the relationship between returns of global phishing crime numbers and Ethereum's key financial metrics (including total number of transactions, average price per transaction, average transaction quantity, aggregate quantity of tokens traded, conditional on the states of the respective markets). The results indicate a significant positive relationship between changes in phishing crime numbers and several of Ethereum's financial metrics, in particular, with large increases in the total number of transactions, average transaction price, and transaction quantity. The findings provide initial insights into the connection between cryptocurrency usage and global phishing crimes, which can aid crime agencies and regulators in designing actionable strategies.

C1178: **Stylized facts of metaverse non-fungible tokens**

Presenter: **Stephen Chan**, American University of Sharjah, United Arab Emirates

Co-authors: Yuanyuan Zhang, Jeffrey Chu

Non-fungible tokens (NFTs) within the metaverse represent a rapidly emerging sector in the digital asset space. The aim is to provide a comprehensive review of the metaverse's history and an analysis of the stylized facts of five metaverse NFTs. Market efficiency, volatility clustering, leverage effects, and the return-volume relationship are examined. Key findings show that all NFT returns exhibit kurtosis values significantly exceeding the standard value of three, indicating more peaked and heavier-tailed distributions than a normal distribution. The Hurst exponent fluctuates between 0.3 and 0.8, indicating relative inefficiency in log returns with varying degrees of trend reinforcement and anti-persistence. The GARCH(1,1) model confirms the presence of volatility clustering, with high persistence of volatility shocks over time, and most NFT returns exhibit a negative leverage effect, where negative returns decrease volatility. These findings provide critical insights for investors, content creators, and policymakers, emphasizing the need for innovative strategies and regulatory considerations in this evolving ecosystem.

C1532: **Bayesian causal inference in the presence of endogenous selection into treatment and spillovers**

Presenter: **Duong Trinh**, University of Graz, Austria

A new approach is introduced that harnesses network or spatial data to identify and estimate direct and indirect causal effects in the presence of selection-on-unobservables and spillovers. The proposed framework nests the Generalised Roy model to explicitly account for endogenous selection into treatment and goes beyond to capture spillovers through exposure mapping to neighbours' treatment. This allows for heterogeneous effects across individuals and enables exploration of various policy-relevant treatment effects. We develop Bayesian estimators based on data augmentation methods, offering efficient computation and proper uncertainty quantification, which is supported by simulation experiments. We apply the method to evaluate the Opportunity Zones (OZ) program, which aims to stimulate economic growth in distressed U.S. census tracts through tax incentives. The results show both direct and indirect positive impacts on housing unit growth in designated Qualified Opportunity Zones (QOZs), but unselected tracts (non-QOZs) experience no beneficial spillovers, remaining at a disadvantage. Moreover, the model predicts that offering investment tax credits to non-QOZs would lead to negative outcomes, making the program's expansion to these areas ineffective.

CC376 Room BCB 309 APPLIED MACHINE LEARNING

Chair: Bettina Gruen

C1200: **Elucidating the temporal dimension of the genotype by environment interaction effect with transformers**

Presenter: **Patrick McMillan**, University of Guelph, Canada

Co-authors: Zeny Feng, Lewis Lukens

Genomic selection (GS) is a powerful tool to predict the phenotype of an organism based solely on its genotype. GS has quickly become integral to large-scale crop breeding programs, as the ability to fit genome prediction (GP) models using high-density markers can accelerate the rate of genetic gain. A limiting factor in the application of GS to crop breeding is the genotype-by-environment interaction effect. This effect confounds the ability of GP models to identify elite lines in a breeding scheme, slowing genetic gain. The aim is to present a novel framework to elucidate the temporal dimension of the genotype-by-environment interaction effect through the use of weekly climate observations, as well as static soil and landscape factors, along with directly incorporating single-nucleotide polymorphism (SNP) data in a modified temporal fusion transformer (TFT) model developed for the prediction of previously unseen crop varieties to emulate a breeding scheme. It is demonstrated that the TFT model is not only able to significantly outperform linear models commonly used in the field, but it can also provide previously unseen insight into the effects of individual climate and SNP variables across the length of a growing season.

C1057: **Probabilistic deep learning for forecasting influenza hospitalizations**

Presenter: **Itai Dattner**, University of Haifa, Israel

Co-authors: David Shulman, Rami Yaari, Jeffrey Shaman

Accurate forecasting of hospitalizations due to influenza is essential for effective public health planning and response. In the U.S., the CDC's FluSight Forecast Hub collects and compares real-time forecasts of influenza-associated hospitalizations from various modeling teams. However, many models face trade-offs between predictive accuracy and interpretability. A deep learning framework is proposed that generates probabilistic forecasts of influenza-related hospitalizations across U.S. states by parameterizing a negative binomial distribution through a neural network. The approach captures both spatial and temporal variability while providing full predictive distributions. Model explainability and theoretical guarantees are discussed for the proposed methodology. Preliminary results indicate that the approach achieves state-of-the-art forecasting accuracy, supporting data-driven public health decision-making.

C1230: **Blockwise boosted inflation: Non-linear determinants of inflation using machine learning**

Presenter: **Galina Potjagailo**, Bank of England, United Kingdom

Co-authors: Marcus Buckmann, Philip Schnattinger

The aim is to propose the Blockwise Boosted Inflation Model (BBIM), a boosted tree framework that decomposes inflation dynamics into predictive components aligned with an open-economy hybrid Phillips curve. Demand and supply contributions are identified by imposing monotonicity constraints, ensuring theory-consistent links between inflation and key indicators. Applied to monthly UK CPI inflation, the model shows that the recent surge has been driven mainly by global supply shocks transmitted through supply chains. An L-shaped Phillips curve relationship between inflation and labour market tightness is also uncovered, with tight labour markets amplifying recent inflationary pressures. By contrast, earlier episodes saw non-linearities more strongly tied to broader slack, particularly during recessions. The model further accounts for trend shifts informed by inflation expectations. Short-term household expectations have recently displayed persistent non-linear effects, temporarily raising trend inflation and prolonging inflationary pressures, while longer-term expectations remain anchored. Out-of-sample, the BBIM delivers competitive forecasting performance relative to linear benchmarks and unstructured machine learning methods. The approach provides a flexible yet interpretable framework that combines economic structure with machine learning for policy-relevant analysis of inflation dynamics.

C0233: **Optimization of electricity demand forecasting using machine learning ensemble methods for Israel's energy grid**

Presenter: **Moshe Kelner**, University of Haifa and Noga - Israel System Operator, Israel

Accurate electricity demand forecasting is critical for Israel's isolated energy grid, producing 80 billion kWh annually from domestic sources. Daily generation unit scheduling relies on demand forecasts, making accuracy essential for economic and environmental optimization. Forecast errors have significant consequences. Underestimation forces expensive backup unit activation, while overestimation causes unnecessary generation costs. An ensemble system is developed, combining five machine learning algorithms to predict daily electricity demand and peak demand for three-day horizons. The approach incorporates 1,200 meteorological features, including hourly temperatures across regions, humidity, and historical patterns. Feature selection used statistical significance testing, with training on 2000-2022 data and 2023 testing. Research addresses three methodological

challenges - establishing performance-based ensemble weighting criteria, improving rare extreme event prediction despite limited examples, and developing adaptive algorithms for consistent performance during seasonal transitions with high weather variability. Model performance was evaluated using MAPE between forecasted and actual values. The ensemble demonstrates improved accuracy over individual models, directly supporting energy sector decisions through optimized scheduling, cleaner source prioritization, reduced costs, and enhanced grid reliability.

CC347 Room BCB 405 BAYESIAN METHODS

Chair: Alejandro Murua

C1266: Bayesian estimation of multiple change points in factor copula models

Presenter: **Marvin Borsch**, University of Cologne, Germany

Co-authors: Roman Liesenfeld, Dominik Wied

The aim is to introduce a methodology for detecting multiple change points in the loadings of a factor copula model, resulting in shifts in the underlying correlation structure. By leveraging a copula framework, the marginal distributions are modeled independently of their dependence structure. A Bayesian procedure is developed to jointly estimate the factor loadings and the locations of the change points, assuming a fixed number of structural changes. Posterior samples are obtained via an independent Metropolis-Hastings within Gibbs sampler. In this approach, the conditional posterior for the loadings depends solely on the likelihood within each segment and its prior, while each change points conditional posterior is determined only by the likelihood of the data between adjacent change points and its prior. The proposed approach is evaluated through Monte Carlo simulations under various change scenarios, and its practical relevance is demonstrated through an empirical application to the correlation structure of EURO STOXX 50 companies across different sectors during the recent Covid-19 pandemic.

C1404: Hierarchical Bayesian calibration with Bayesian committee machine

Presenter: **Sebastian Heinekamp**, Paul Scherrer Institut, Switzerland

Co-authors: David Higdon

Calibrating computer simulation code to experimental observations is an inherent task in any field using simulations to guide and inform experiments. Motivated by the goal to understand parameter uncertainties in particle accelerator experiments, hierarchical Bayesian calibration is used. Certain inputs (such as the beam injection amplitude) need to be calibrated separately for each experiment. The approach follows the calibration model suggested by a prior study, while a hierarchical description of the collection of parameters inferred for the individual experiments is added. The implementation makes use of a large number of simulations, leading to computational challenges for posterior exploration via Markov chain Monte Carlo (MCMC). To overcome this, the Bayesian committee machine (BCM) approximation is used for large GPs, allowing for leveraging high-performance computing infrastructure to speed up the resulting MCMC. The BCM offers problem and dataset-agnostic speed up and parallelization. Using the No-U-turn sampling algorithm allows us to leverage Julia's automatic differentiation capabilities. The improvements in runtime are demonstrated on analytic examples and simulation data for the Argonne Wakefield Accelerator.

C1335: Sampling from density power divergence-based generalized posterior distribution via stochastic optimization

Presenter: **Naruki Sonobe**, Tokyo University of Science, Japan

Co-authors: Tomotaka Momozaki, Tomoyuki Nakagawa

Robust Bayesian inference using density power divergence (DPD) has emerged as a promising approach for handling outliers in statistical estimation. While the DPD-based posterior offers theoretical guarantees of robustness, its practical application to general parametric models is computationally challenging due to an analytically intractable integral term in its formulation. These challenges are specifically pronounced in high-dimensional settings, where traditional numerical integration methods are inadequate and computationally expensive. The aim is to propose a new approximate sampling methodology that addresses these limitations by integrating the loss-likelihood bootstrap with a stochastic gradient descent algorithm specifically designed for DPD-based estimation. The approach enables efficient and scalable sampling from DPD-based posteriors for a broad class of parametric models, including those with intractable integrals. It is further extended to accommodate generalized linear models. Through simulations, it is demonstrated that the method generates samples that accurately approximate the target posterior while offering superior computational scalability. The results confirm that this framework provides a practical and efficient tool for applying robust Bayesian inference to complex, high-dimensional data.

C1333: Empirical Bayes 1-bit matrix completion

Presenter: **Takeru Matsuda**, University of Tokyo & RIKEN Center for Brain Science, Japan

The problem of predicting unobserved entries of a binary data matrix from observed entries is called the 1-bit matrix completion. An empirical Bayes method is developed for 1-bit matrix completion that utilizes a low-rank structure like the multidimensional item response theory. The proposed method is motivated by an empirical Bayes estimator of a normal mean matrix of a prior study, which is a matrix generalization of the James–Stein estimator and shrinks the singular values towards zero. Simulation results and application to real data are presented.

Monday 15.12.2025

16:50 - 18:30

Parallel Session O – CFE-CMStatistics 2025

CI009 Room BCB 206 CFE SPECIAL INVITED SESSION: A TRIBUTE TO H. PESARAN III**Chair: Natalia Bailey****C0164: High-dimensional dynamic factor models: A selective survey***Presenter:* **Manfred Deistler**, Vienna University of Technology, Austria*Co-authors:* Marco Lippi, Brian Anderson

High-dimensional dynamic factor models are presented in detail: The main assumptions and their motivation, main results, illustrations by means of elementary examples. In particular, the role of singular ARMA models in the theory and applications of high-dimensional dynamic factor models is discussed. The emphasis is on model classes and their structure theory, rather than on estimation in the narrow sense. The survey is not comprehensive. Its aim is to point out promising lines of research and applications that have not yet been sufficiently developed.

C0165: Generalized impulse responses, multi-horizon projections, and causal mediation analysis in macroeconomics and finance*Presenter:* **Jean-Marie Dufour**, McGill University, Canada

A novel concept of impulse response decomposition is introduced to disentangle the dynamic contributions of the mediator variables in the transmission of structural shocks. The decomposition is justified by drawing on causal mediation analysis and demonstrating its equivalence to the average mediation effect. The result establishes a formal link between Sims and Granger causality. Sims causality captures the total effect, while Granger causality corresponds to the mediation effect. A dynamic mediation index that quantifies the evolving role of mediator variables in shock propagation is constructed. Applying the framework to studies of the transmission channels of U.S. monetary policy, it is found that investor sentiment explains approximately 60 percent of the peak output response in three months following a contractionary policy shock, while expected default risk contributes negligibly across all horizons.

C0713: Financial statements and macroeconomic dynamics*Presenter:* **Allan Timmermann**, UCSD, United States*Co-authors:* Davide Pettenuzzo, Riccardo Sabbatucci

What do companies' filings reveal about the state of the macro economy, and do specific accounting variables contain particularly relevant information? To address these questions, the lead-lag patterns of more than twenty accounting variables are analyzed in relation to aggregate economic activity. New daily corporate account business activity indices are developed, which aggregate firm-level accounting information while controlling for shifts in the composition of announcers and reducing firm-specific noise. The new indices show that firm liquidity becomes significantly lower while corporate debt grows significantly faster several months prior to recessions, and thus can be used as leading indicators. Conversely, operations, earnings, and profitability measures tend to be significantly lower after recessions, suggesting they are mostly lagging, pro-cyclical indicators of economic activity.

CO071 Room BCB G07 HiTEC: ADVANCES IN MATRIX TIME SERIES IN ECONOMETRICS**Chair: Ivan Ricardo****C1093: Detecting cointegrating relations in non-stationary matrix-valued time series***Presenter:* **Ivan Ricardo**, Maastricht University, Netherlands

The aim is to propose a matrix error correction model to identify cointegration relations in matrix-valued time series. Separate cointegrating relations are hereby allowed along the rows and columns of the matrix-valued time series, and information criteria are used to select the cointegration ranks. Through Monte Carlo simulations and a macroeconomic application, it is demonstrated that the approach provides a reliable estimation of the number of cointegrating relationships.

C1101: Bayesian dynamic factor models for high-dimensional matrix-valued time series*Presenter:* **Wei Zhang**, Johns Hopkins University, United States

High-dimensional matrix-valued time series are of significant interest in economics and finance, with prominent examples including cross-regional macroeconomic panels and firms' financial data panels. The aim is to introduce a class of Bayesian matrix dynamic factor models that utilize matrix structures to identify more interpretable factor patterns and factor impacts. The model accommodates time-varying volatility, adjusts for outliers, and allows cross-sectional correlations in the idiosyncratic components. For model comparison, an importance-sampling estimator is employed based on the cross-entropy method to inform decisions regarding: (1) the optimal dimension of the factor matrix; (2) the appropriate factor structure—whether vector-valued or matrix-valued; and (3) the suitability of an approximate versus exact factor model. Through a series of Monte Carlo experiments, the properties of the factor estimates and the performance of the marginal likelihood estimator in correctly identifying the true model are shown. Applying the model to a macroeconomic dataset and a financial dataset, its ability is demonstrated in unveiling interesting features within matrix-valued time series.

C1181: Cointegrated models for matrix valued time-series*Presenter:* **Emanuele Lopetuso**, Università di Padova, Italy*Co-authors:* Massimiliano Caporin

Traditional econometric analyses represent observations as vectors despite the inherent complexity of empirical data structures. When data are organized along dual classification dimensions, a matrix representation provides a more natural and interpretable framework. Building on recent advances in matrix autoregressive (MAR) modeling, the aim is to introduce a novel error correction representation tailored for matrix-structured data. Through comparative analysis with existing methodologies, we demonstrate two critical advancements. First, the proposed model preserves the interpretative foundations of conventional cointegration analysis, with coefficients that explicitly capture dynamics rooted in adjustment toward steady-state positions. Second, in contrast to previous formulations, the error correction framework allows for an equivalent matrix autoregressive representation, preserving the fundamental structure of the data in both specifications. This ensures that the matrix representation reflects an intrinsic characteristic of the data.

C1159: Inference in matrix-valued time series with common stochastic trends and multifactor error structure*Presenter:* **Greta Goracci**, Free University of Bozen/Bolzano, Italy*Co-authors:* Rong Chen, Lorenzo Trapani, Simone Giannerini

The aim is to develop an estimation methodology for a factor model for high-dimensional matrix-valued time series, where common stochastic trends and common stationary factors can be present. The focus, in particular, is on the estimation of (row and column) loading spaces, of the common stochastic trends and of the common stationary factors, and the row and column ranks thereof. In a set of (negative) preliminary results, it is shown that a projection-based technique fails to improve the rates of convergence compared to a flattened estimation technique, which does not take into account the matrix nature of the data. Hence, a three step algorithm is developed where: (i) the data is first projected onto the orthogonal complement to the (row and column) loadings of the common stochastic trends; (ii) such trend free data is subsequently used to estimate the stationary common component; (iii) the estimated common stationary component is removed from the data, and re-estimated, using a projection-based estimator, the row and column common stochastic trends and their loadings. It is shown that this estimator succeeds in refining the rates of convergence of the initial, flattened estimator. As a byproduct, consistent eigenvalue ratio-based estimators are developed for the number of stationary and nonstationary common factors.

CO183 Room BCB G08 EMPIRICAL MACRO-FINANCE**Chair: Aikaterini Karadimitropoulou****C0371: Sweet, crude, and golden: Superior commodity price forecasting with penalized linear methods****Presenter:** Juan Laborda, Universidad Politecnica Madrid, Spain**Co-authors:** James Chen, Charalampos Agiropoulos

A structured framework is proposed for commodity price forecasting, balancing interpretability, predictive accuracy, and computational efficiency. A four-stage progression is constructed: (1) A CAPM+4 model extending the capital asset pricing model beyond volatility to skewness and kurtosis; (2) penalized methods of linear regression (Ridge, Lasso, ElasticNet) to mitigate multicollinearity and overfitting; (3) Gaussian process regression (GPR) for quantifying uncertainty; and (4) machine learning ensembles (random forests, extra trees, XGBoost). Prices for five commodities are forecasted: Brent, natural gas, copper, gold, and sugar. Penalized linear methods, notably a Bayesian implementation of the L2 penalty, outperform both OLS and machine ensembles, while maintaining computational tractability and economic interpretability. Additional innovations include RAFF (regularization-adjusted factor estimation), a novel diagnostic derived from Bayesian ridge hyperparameters that identifies periods of market stress. Furthermore, stacked generalization produces a meta-forecast. Findings challenge the presumption that complex ML models dominate financial forecasting. Penalized linear methods hit the "sweet spot", combining performance, speed, and transparency. The framework provides clear criteria for model selection based on interpretability, data heterogeneity, and computational constraints.

C0387: Human capital development measurement in a digital world**Presenter:** Alexandros Bechlioulis, University of Piraeus, Greece**Co-authors:** Claire Economidou, Nikolas Topaloglou

Digital technologies have moved beyond the margins to become transformative drivers reshaping the very fabric of global human capital development. The purpose is to challenge the adequacy of the United Nations Human Development Index (HDI) - a cornerstone metric in development studies - by proposing a fundamental expansion: The inclusion of digitalization as a critical fourth dimension. Two innovative methodologies are employed within the stochastic dominance framework: A stochastic bounding analysis to construct a digital performance index and a stochastic spanning test to evaluate welfare implications. Results reveal a strikingly different global development narrative. When human capital development is assessed through the proposed digitally-augmented HDI measure, traditional country rankings, based on the standard HDI, and developmental pathways are profoundly reshaped, exposing digitalization not only as a potent catalyst for advancement but also as a source of persistent and widening digital inequalities.

C0639: Commodity price uncertainty comovement: Does it matter for global economic growth?**Presenter:** Aikaterini Karadimitropoulou, University of Piraeus, Greece**Co-authors:** Laurent Ferrara, Athanassios Triantafyllou

Using a dynamic factor model, a global commodity uncertainty factor is estimated, which captures comovement in volatilities of major agricultural, metals, and energy commodities. Then, impulse response functions computed via a structural VAR model show that an increase in this global shock results in a substantial drop in investment and trade, for both emerging and advanced economies. Last, the methodology disentangles "good" and "bad" macroeconomic effects of oil price uncertainty: An oil price uncertainty shock common to all commodities leads to recessionary effects. However, positive short-run effects are observed when this shock is only specific to the oil market.

C0727: The dynamics of velocity of money and money demand in the United States**Presenter:** Krzysztof Beck, Lazarski University, Poland

In the first part of the investigation, we analyze the dynamics of velocity of money in the United States by examining M0, M1, M2, and M3 monetary aggregates within the context of the quantity theory of money. The long-run relationship implied by the quantity theory of money is estimated, allowing for time-varying velocity of money and short-run dynamics, using a Bayesian error-correction mechanism model. These changes are examined within a Markov switching VEC (MS-VEC) model and a specification that permits the cointegrating space to evolve over time in a manner comparable to the random walk variation (TVP-VEC). Using the Frobenius matrix norm, the distance between the obtained cointegration spaces and those spanned by the vectors is measured based on the money volatility reported in the Federal Reserve of Saint Louis Database. In the second part of the investigation, it is estimated how demand for money depends on interest rates, income, inflation, and expected inflation, as well as the uncertainty underlying those determinants. It is found that the volatility of money behaves very differently from the basic estimates depicted in the Federal Reserve Bank of St. Louis Database.

C0776: Fiscal rules, debt surprises, and stock flow adjustments**Presenter:** Marcos Poplawski Ribeiro, International Monetary Fund, United States

The purpose is to analyze whether fiscal rules reduce forecast errors on public debt and its components, in particular, stock-flow adjustments. Using multi-year forecasts from annual vintages of the IMF's World Economic Outlook database for 176 countries during 1996-2023, it is shown that countries with fiscal rules in place have a more accurate projection of stock-flow adjustments. Moreover, governments that comply with their fiscal rules have a more predictable public debt trajectory. These findings indicate the importance of the public financial management institutions related to fiscal rules, such as those reporting and monitoring compliance with the rules.

CO057 Room BCB G09 TOPICS IN ECONOMETRICS**Chair: Johan Lyhagen****C0225: Design-based inference under random potential outcomes via Riesz representation****Presenter:** Yukai Yang, Uppsala University, Sweden

A design-based framework is introduced for causal inference that accommodates random potential outcomes, thereby extending the classical Neyman-Rubin model in which outcomes are treated as fixed. Each unit's potential outcome is modeled as a structural mapping $\bar{y}_i(z, \omega)$, where z denotes the treatment assignment and ω represents latent outcome-level randomness. Inspired by recent connections between design-based inference and the Riesz representation theorem, potential outcomes are embedded in a Hilbert space and define treatment effects as linear functionals, yielding estimators constructed via their Riesz representers. This approach preserves the core identification logic of randomized assignment while enabling valid inference under stochastic outcome variation. Large-sample properties are established under local dependence, and consistent variance estimators that remain valid under weaker structural assumptions are developed, including partially known dependence. A simulation study illustrates the robustness and finite-sample behavior of the estimators. Overall, the framework unifies design-based reasoning with stochastic outcome modeling, broadening the scope of causal inference in complex experimental settings.

C0400: Levels or directions: Forecasting forward freight agreements**Presenter:** Jonas Andersson, Norwegian School of Economics, Norway**Co-authors:** Lisa Maria Assmann, Roar Adland

The aim is to investigate how to best forecast the forward market for maritime freight. Complementing existing studies that are mostly based on statistical evaluation measures, the focus is on the economic consequences of predictions and mimicking a trader's decision in the freight forward market. This leads to a comparison of directional and point forecasting methods, concluding that predictions based on past directions are significantly better economically. With the unrealistic assumption that trading signals based on closing prices can be executed without delay,

significant profits are found. If a more onerous assumption of one day price slippage is implemented, prompted by signs of illiquidity in the FFA market, resulting profits are diminished, challenging the robustness of the conclusion of an inefficient FFA market in practice.

C0624: **Nonlinear vector autoregressive models and unit roots**

Presenter: **Rickard Sandberg**, Stockholm School of Economics, Sweden

The purpose is to investigate the challenges of testing for linearity in the context of univariate and multivariate smooth transition autoregressive (STAR and VSTAR) models. A key focus is the impact of unit root processes on the distribution of linearity test statistics, which can lead to substantial size distortions. These distortions pose a risk of spurious rejections of linearity, potentially resulting in model misspecification and unreliable inference. For example, a test nominally set at the 5% significance level may exhibit empirical rejection rates as high as 17.1% in univariate STAR models and 22.6% in bivariate VSTAR models. Critically, these distortions are shown to magnify with the dimensionality of the system; in a VSTAR model with four variables, the empirical size reaches 46.1%. A novel contribution is the derivation of the asymptotic distributions of linearity tests under a unit root assumption in multivariate systems - an extension not previously addressed in the literature. An empirical application using U.S. inflation and Federal Funds Rate data illustrates the importance of accounting for unit roots prior to testing for nonlinearity. The findings emphasize the need for adjusted critical values and robust testing strategies when working with persistent and high-dimensional macroeconomic time series.

C0822: **A review of risk aversion coefficient estimation for quadratic and exponential utility functions: Empirical evidence**

Presenter: **Sachi Dilami Ilangasekara**, Linnaeus University, Sweden

Co-authors: Peter Karlsson, Stanislas Muhinyuza

The aim is to present a comprehensive review of the methodology used to estimate the risk aversion coefficient, focusing on quadratic and exponential utility functions. It examines the theoretical foundations and practical applications of these utility frameworks under various distributional assumptions of asset returns. Furthermore, it compares different estimators of the risk aversion coefficient, analyzing and assessing their statistical properties under the distributional assumptions considered. A structured comparison is provided to support the selection of appropriate estimators for the risk aversion coefficient, depending on the underlying distribution assumptions and utility framework.

CO229 Room Virtual R01 BAYESIAN METHODS FOR COMPOSITIONAL AND PROPORTIONAL DATA

Chair: Matthew Heiner

C0333: **Scalable and robust regression models for continuous proportional data**

Presenter: **Changwoo Lee**, Duke University, United States

Co-authors: Benjamin Dahl, Otso Ovaskainen, David Dunson

Beta regression is used routinely for continuous proportional data, but it often encounters practical issues such as a lack of robustness of regression parameter estimates to misspecification of the beta distribution. An improved class of generalized linear models is developed, starting with the continuous binomial (cobin) distribution and further extending to dispersion mixtures of cobin distributions (micobin). The proposed cobin regression and micobin regression models have attractive robustness, computation, and flexibility properties. A key innovation is the Kolmogorov-Gamma data augmentation scheme, which facilitates Gibbs sampling for Bayesian computation, including in hierarchical cases involving nested, longitudinal, or spatial data. Robustness, ability to handle responses exactly at the boundary (0 or 1), and computational efficiency relative to beta regression are demonstrated in simulation experiments and through analysis of the benthic macroinvertebrate multimetric index of US lakes using lake watershed covariates.

C0697: **A Bayesian semiparametric mixture model for clustering zero-inflated microbiome data**

Presenter: **Matthew Koslovsky**, Colorado State University, United States

Microbiome research has immense potential for unlocking insights into human health and disease. A common goal in human microbiome research is identifying subgroups of individuals with similar microbial composition that may be linked to specific health states or environmental exposures. However, existing clustering methods are often not equipped to accommodate the complex structure of microbiome data and typically make limiting assumptions regarding the number of clusters in the data, which can bias inference. Designed for zero-inflated multivariate compositional count data collected in microbiome research, a novel Bayesian semiparametric mixture modeling framework is proposed that simultaneously learns the number of clusters in the data while performing cluster allocation. In simulation, the clustering performance of the method is demonstrated compared to distance- and model-based alternatives, and the importance of accommodating zero-inflation when present in the data. The model is then applied to identify clusters in microbiome data collected in a study designed to investigate the relation between gut microbial composition and enteric diarrheal disease.

C0723: **Shrinkage on the simplex: A Bayesian framework for quantifying sparsity and dependence in compositional data**

Presenter: **Jyotishka Datta**, Virginia Polytechnic Institute and State University, United States

Co-authors: Matthew Heiner, David Dunson, Otso Ovaskainen

Sparse signal recovery remains a central challenge in large-scale data analysis. Over the past decade, global-local shrinkage priors have emerged as the Bayesian gold standard for sparse inference and a wide range of nonlinear problems. Yet, discrete compositional data, routinely encountered in fields like microbiomics, pose unique difficulties: The Dirichlet distribution cannot adapt to arbitrary levels of sparsity in high-dimensional probability vectors. A new shrinkage prior is introduced on the simplex, specifically designed to scale to problems with many categories while flexibly capturing both sparsity and dependence among components. Theoretical properties that guarantee adaptive behavior are presented, an efficient posterior sampling scheme is outlined, and through simulations and an application to microbiome data, it is demonstrated that the approach outperforms both standard Dirichlet models and existing alternatives.

C0872: **Partially geometric stick-breaking processes**

Presenter: **Matthew Heiner**, Brigham Young University, United States

Co-authors: Gilbert Fellingham, Alejandro Jara, Garritt Page

Empirical distributions for count data often exhibit idiosyncrasies among low values, such as zero inflation and stable tail behavior. A flexible model is proposed for counts that combines nonparametric estimation of early probabilities with fixed decay after a single, unknown change point. A stick-breaking construction is used with variables that are beta distributed before the change point and share a single value thereafter, inducing a geometric tail. The resulting process model parsimoniously balances needed flexibility where data are abundant, with a parametric representation where data are naturally sparse. The construction admits a collapsed posterior distribution for the change point, avoiding transdimensional MCMC. It is illustrated by modeling rally lengths in men's professional tennis, where the change point may indirectly measure the effect of server advantage. Finally, the new process is used to construct a countably infinite mixture model, extending the geometric stick-breaking process of Fuentes-Garcia, Mena, and Walker, and demonstrating its effectiveness for density estimation. The random change point accommodates a variety of behaviors, ranging from a parsimonious approximation of Poisson-Dirichlet process mixtures to mixtures of finite mixtures.

CO240 Room BCB 208 SYNTHETIC DATA: GENERATION AND VALIDATION METHODS

Chair: Nora Amama Ben Hassun

C0236: **Mind the gap: From synthetic data to regulatory confidence in healthcare AI**

Presenter: **Vibeke Binz Vallevik**, University of Oslo, Norway

The growing use of synthetic data in healthcare AI demands validation strategies that ensure both technical integrity and regulatory compliance.

Current validation practices largely emphasize statistical similarity to real data, falling short of the requirements outlined by the EU's MDR, IVDR, and AI Act. Moreover, the GDPR's ambiguous definition of personal data creates uncertainty about the legal status of synthetic health data, complicating its use in development and approval pipelines. At the same time, growing regulatory openness to synthetic data in product development and approval underscores the need for multidimensional validation approaches. The EU-funded SYNTHIA project is addressing this gap by developing a quality assurance framework aligned with EU regulatory principles of risk management, transparency, and clinical relevance. The framework extends beyond conventional metrics to incorporate fairness, privacy, and environmental impact. A case study using membership inference attacks on a synthetic cancer dataset demonstrates practical methods to assess residual privacy risks. In parallel, legal analysis based on fundamental rights helps clarify when synthetic data might still fall under GDPR, guiding developers on compliance. Together, these approaches support a structured path for responsible and lawful integration of synthetic data into healthcare AI.

C0425: Differentially private synthetic data without training

Presenter: **Zinan Lin**, Microsoft Research, United States

Generating differentially private (DP) synthetic data that closely resembles original data while preserving user privacy is a scalable solution to address privacy concerns in today's data-driven world. Private evolution (PE) is introduced, a new training-free framework for DP synthetic data generation, which contrasts with existing approaches that rely on training DP generative models. PE treats foundation models as blackboxes and only utilizes their inference APIs. It is demonstrated that across both images and text, PE: (1) could match or even outperform prior state-of-the-art (SoTA) methods in the fidelity-privacy trade-off without any model training in some cases; (2) enables the use of advanced open-source models (e.g., Mixtral) and API-based models (e.g., GPT-3.5), where previous SoTA approaches are inapplicable; and (3) is more computationally efficient than prior SoTA methods. Additionally, recent extensions of PE are discussed, including the integration of non-neural-network data synthesis tools, fusion of knowledge from multiple models for DP data synthesis, and applications in federated learning. The hope is that PE unlocks the full potential of foundation models and other data synthesis tools in privacy-preserving machine learning and accelerates the adoption of DP synthetic data across industries.

C0517: Synthetic tabular data detection in the wild

Presenter: **Gaspard Charbel Novixi Kindji**, Orange Innovation, France

The rapid progress of generative models offers remarkable capabilities, but also raises data integrity concerns, especially in distinguishing authentic from synthetic data. This concern has gained significant attention in the realms of image and text. However, for data types such as tabular data, the landscape of generative models is getting richer, but little attention is paid to detection techniques. Detecting synthetic tabular data is uniquely difficult due to its heterogeneous and variable structure, with the main difficulty lying in data representation rather than the classifier itself. The challenge of detecting synthetic tabular data in real-world scenarios, where detectors must generalize to unseen table formats, is addressed. A novel datum-wise transformer architecture is introduced, designed to operate effectively across arbitrary tabular structures, which encodes all features as text with an independent embedding for each feature. The method achieves an AUC of 0.67 compared to 0.60 for the best competitor. This result is later improved to 0.69 with domain adaptation techniques. The result is the reliable, scalable detection of synthetic tabular data, which can be extended to other mainstream predictive tasks involving tabular data.

C0843: Enhanced validation of tabular synthetic data: Assessing propensity score resemblance metrics

Presenter: **Nora Amama Ben Hassun**, Universitat Politècnica de Catalunya, BarcelonaTech (UPC), Spain

Co-authors: Daniel Fernandez, Jordi Cortes Martinez

Rigorous assessment of validation metrics is a prerequisite for a unified, variable-class-aware framework for tabular synthetic data. In particular, the assessment of resemblance by multivariate metrics quantifies OD and SD similarity, guides synthesizer refinement, and standardizes method comparison. A simulation study, followed by a real data case study, was implemented in R using the synthpop package. Synthetic datasets are generated over different sample sizes (n) and number of variables (p) under the hypothesis of OD and SD coming from the same population. Null distributions of three propensity score-based metrics were derived to quantify type I error. To evaluate statistical power, alternative scenarios introduced controlled shifts in means, variances, intervariable correlations, and distributional symmetry. Propensity scores were estimated via logistic regression with a train and test split to prevent classifier overfitting. All three metrics controlled type I error; scenarios are also delineated where each metric fails. The metric is identified with superior statistical power across alternative scenarios. The findings suggest that certain metrics may be employed to validate the resemblance of SD from a multivariate perspective in the context of numerical variables. Further research is required in order to explore the expansion for different classes of variables.

C0330 Room BCB 209 LARGE DATA ANALYSIS IN CANCER GENOMICS AND STATISTICS FOR CANCER Chair: Damianos Michaelides

C0917: Bayesian clustering of prostate cancer patients with simultaneous feature selection of metabolites

Presenter: **Inga Huld Armann**, Imperial College London, United Kingdom

Co-authors: Elizabeth Bancroft, Zsolt Kote-Jarai, Ros Eeles, Ioanna Papatouma, Marina Evangelou

Prostate cancer (PrCa) is the most common cancer in men in the UK, with 55,000 cases and 12,000 deaths per year. Leveraging data obtained from accessible screening methods, such as blood tests, would enhance early diagnosis and improve patient care. Metabolomics extracted from biofluids such as blood offer a promising and easily accessible source of data and have been studied alongside other 'omics' datasets to advance cancer research. Metabolomics data from two case-control studies on PrCa, namely, PROFILE and IMPACT, were analyzed. Building on the success of previous analyses of metabolomics data for PrCa, clustering techniques are employed to identify potential cancer subtypes, particularly those linked to disease aggressiveness. Motivated by the importance of biomarker identification, Bayesian nonparametric clustering via Dirichlet process mixture models is proposed, with an extension allowing for simultaneous feature selection. The extension aims to identify features contributing to the clustering of observations. This is an important insight into potential biomarkers. The posterior is approximated through variational inference, offering a computationally efficient alternative to traditional Markov Chain Monte Carlo methods. The effectiveness of the proposed method is evaluated on both real and simulated data. Moreover, the method is applicable to other 'omics' datasets, including genomics.

C0416: Large scale case-control analyses using LR modelling improves rare variant classification in breast cancer risk genes

Presenter: **Damianos Michaelides**, The Cyprus Institute of Neurology and Genetics, Cyprus

A large number of rare variants in breast cancer susceptibility genes remain as variants of uncertain significance (VUS). The impact of hundreds of VUS is quantified through evidence derived from analysis of over 15K variants in PALB2, TP53, CHEK2, and ATM. The statistical methodology employed benefits from the largest breast cancer case-control dataset to date (>60K cases, >250K controls), using sequencing data. The cCLR method is used, which models the likelihood of observing a variant in cases versus controls, incorporating survival information and known gene-level penetrance. However, the method assumes a uniform risk across all variants in a gene, ignoring that individual variants may confer stronger or weaker risks, hence, limiting its ability to detect risk heterogeneity. This is addressed by a refined approach using a dynamic scaling parameter that calibrates relative risk per-variant. This allows the model to evaluate how different risk magnitudes fit the data for each variant. The refined method is useful in identifying variants that carry a higher or lower risk than the gene's average pathogenic variants. LRs were computed across a range of scaling values, representing a spectrum of relative risks. Pathogenic evidence was based on the maximum LR. Benign evidence was based on the LR at the gene-level risk. The refined approach provides evidence for 368 rare unclassified variants, a 15% increase over the evidence identified by the standard method.

C1136: Spatial transcriptomics reveals gene programs of plasma cell rich lymphomyeloid aggregates in ovarian cancer*Presenter:* **Shreena Nisha Kalaria**, Deeley Research Centre at BC Cancer Victoria, Canada*Co-authors:* Celine M Laumont, Farouk Nathoo, Brad Nelson

Plasma cells are emerging as prognostic markers in high-grade serous ovarian cancer (HGSC), yet their spatial organization and transcriptional states within the tumour microenvironment remain poorly characterized. In particular, plasma cell-rich lymphomyeloid aggregates (LMAs) - dense immune cell regions in which plasma cells dominate - are potentially important sites of anti-tumour activity. A modified 10x Genomics Visium workflow is used to obtain spatially resolved gene expression and B-cell receptor (BCR) data from 14 HGSC tumours. Plasma cell-rich LMAs were identified via kernel density estimation, and their transcriptional landscapes were characterized using negative binomial differential expression tests. Three predominant transcriptional programs emerged: a ribosomal-low state, an extracellular matrix remodelling state, and an innate immunity activation state. These likely represent distinct functional states. Preliminary analyses suggest that plasma cell-rich LMAs may exhibit distinct proximities to epithelial structures compared with other immune aggregates, as well as unique combinations of BCR receptors. These relationships will be quantified, and their reproducibility will be assessed in external datasets. Findings demonstrate how integrating spatial statistics with transcriptomic profiling can reveal the heterogeneity of immune microenvironments in HGSC, laying the groundwork for understanding their therapeutic relevance.

C1138: FiXeD: Spatial point process distances for pairing the heavy and light chains of B cell receptors from spatial BCR-seq*Presenter:* **Yimeng Liu**, University of Victoria, Canada

The immune system identifies and mounts a defense against tumor cells through antigen recognition, a process mediated primarily by T cells and B cells. B cells are equipped with heavy and light receptor chains, whereas T cells carry alpha and beta receptor chains. These receptors recognize antigens and orchestrate an immune response, making their accurate pairing essential for understanding tumor-immune interactions. The purpose is to develop an approach based on spatial transcriptomic data and BCR-seq to infer receptor chain pairings. The approach is formulated as a combinatorial optimization problem with an objective function that incorporates the expression matrices of heavy and light chains with their spatial co-expression point patterns derived from barcoded pixels. Spatial information is represented as distances between point patterns computed using an optimal transport algorithm. The estimated receptor heavy and light chain pairs have potential application in the development of targeted immunotherapies. The methodology is evaluated using both simulation studies based on ovarian and breast cancer cohorts, as well as with real data, where the ground truth is known from single-cell sequencing. Substantial improvements in terms of both accuracy and stability are demonstrated over existing state-of-the-art.

CO334 Room BCB 212 NETWORK ECONOMETRICS**Chair: Santiago Pereda-Fernandez****C0600: Peer effects with misspecified peer groups***Presenter:* **Christiarn Rose**, University of Queensland, Australia*Co-authors:* Lizi Yu

The purpose is to consider the identification of peer effects under peer group misspecification. Two leading cases are missing data and peer group uncertainty. Missing data can take the form of some individuals being entirely absent from the data. The researcher does not need to have any information on missing individuals and does not even need to know that they are missing. It is shown that peer effects are nevertheless identifiable under mild restrictions on the probabilities of observing individuals, and a GMM estimator is proposed to estimate the peer effects. In practice, this means that the researcher only needs to have access to an individual-level sample with group identifiers, rather than a sample of groups. Group uncertainty arises when the relevant peer group for the outcome under study is unknown. It is shown that peer effects are nevertheless identifiable if the candidate groups are nested within one another, and a non-linear least squares estimator is proposed. A Monte-Carlo experiment is conducted to demonstrate the identification results and the performance of the proposed estimators, and the method is applied to study peer effects in the career decisions of junior lawyers.

C0637: Regression discontinuity designs under interference*Presenter:* **Tiziano Arduini**, Tor Vergata University of Rome, Italy*Co-authors:* Laura Forastiere, Elena Dal Torrone

The continuity-based framework is extended to regression discontinuity designs (RDDs) to identify and estimate causal effects in the presence of interference when units are connected through a network. In this setting, the assignment to an "effective treatment", which comprises the individual treatment and a summary of the treatment of interfering units (e.g., friends, classmates), is determined by the unit's score and the scores of other interfering units, leading to a multiscore RDD with potentially complex, multidimensional boundaries. These boundaries are characterized, and generalized continuity assumptions are derived to identify the proposed causal estimands, i.e., point and boundary causal effects. Additionally, a distance-based nonparametric estimator is developed, its asymptotic properties are derived under restrictions on the network degree distribution, and a novel variance estimator is introduced that accounts for network correlation. Finally, the methodology is applied to the PROGRESA/Oportunidades dataset to estimate the direct and indirect effects of receiving cash transfers on children's school attendance.

C0846: Partial identification of treatment response under complementarity and substitutability*Presenter:* **Edoardo Rainone**, Banca d'Italia - Dipartimento ECS - Divisione Segreteria, Italy

The purpose is to study partial identification of treatment response when the outcome of an agent is affected heterogeneously by the outcomes and, consequently, the treatment statuses of other agents in the economy. When the sign of interactions is predicted by the theory, a general approach is proposed that allows the use of comparative statics under monotonic treatment response. New theoretical results are provided on how the heterogeneous fixed points theorem can be employed to microfound sharp bounds on the distribution of potential outcomes. It is shown with an empirical application that the method can produce narrow and meaningful bounds for the response to central bank funding of credit to the real economy, under full endogeneity of and heterogeneous interdependence among banks' balance sheet items.

C0952: Inferring social connectedness with endogenously selected groups*Presenter:* **Eric Auerbach**, Northwestern University, United States

Researchers often study social connectedness by first clustering individuals into groups based on the observed network ties. However, standard inference procedures typically ignore the fact that the group definitions are data-driven, leading to biased estimators and invalid inferences. A new method is proposed for inferring social connectedness that remains valid even when the groups are endogenously selected. Theoretical guarantees are first established. Then the approach is evaluated via simulations and in an empirical application.

CO178 Room BCB 213 RECENT ADVANCES IN PANEL TIME SERIES METHODS AND APPLICATIONS**Chair: Peter Pedroni****C1442: Inflationary effects of tariffs and the role of firm level characteristics: A 2-D Bernstein polynomial quantile approach***Presenter:* **Utsav Bahl**, Cambridge University, United Kingdom*Co-authors:* Peter Pedroni

The purpose is to combine macro-level data with a large-scale Compustat panel of firm-level observations to investigate the role of firm-level characteristics in determining the timing and magnitudes of the transmission of tariffs to inflationary responses. In particular, a blended local projection and structurally identified vector auto-regression approach are used to obtain firm-level markup responses to two structurally identified macro-level broad economy-wide tariff shocks: an input cost component and a domestic competition component. The firm-level responses to these

shocks are used as elements in two-dimensional Bernstein polynomial quantile regressions, for which the response horizon represents one of the array dimensions and firm characteristics represent the other array dimension. This framework allows investigating the dynamic consequences of counterfactuals at the firm-level by altering the distribution of firm characteristics to assess how targeted interventions would change aggregate pass-through. Preliminary results suggest that a tariff shock induces a reallocation of sales toward low-markup firms, dampening the aggregate markup response, yet inflationary pressures persist. Importantly, financial structure emerges as a key margin: firms with shorter debt maturities display stronger inflationary pass-through, and the counterfactuals indicate that targeted extensions of debt horizons can directly mitigate tariff-induced inflationary pressures.

C1437: Estimation bias in local projections with heterogeneous panels

Presenter: **Benjamin Alexander**, Federal Reserve Bank of Chicago, United States

Co-authors: Peter Pedroni

Over the past two decades, local projections have become an important tool for estimating dynamics in macroeconomic data. While their properties in the time series context have been well studied, the behavior of local projections on panel data is not fully understood, despite numerous empirical implementations. The focus is on the econometric properties of local projections in a heterogeneous panel context. An analytical expression is derived for the asymptotic bias when heterogeneous dynamics are pooled across units of a stationary panel. This bias is shown to arise from a lagged dependent variable specification, and is shown to depend on the degree of dispersion in dynamics across units. On the basis of this asymptotic expression, finite sample biases with Monte Carlo simulations are documented. A consideration of feasible solutions are concluded with. Findings suggest the importance of appropriately accounting for heterogeneity in panel analysis with local projections.

C1481: Industrial natural gas demand in selected MENA countries. CCE-MG coupled with Autometrics-a machine learning algorithm

Presenter: **Fakhri Hasanov**, KAPSARC, Saudi Arabia

Industrial natural gas demand is examined in selected MENA countries using both panel and country-level analyses. We apply the Common Correlated Effects Mean Group (CCE-MG) estimator and Autometrics, a machine-learning algorithm, under super-saturation. Data for Algeria, Egypt, Bahrain, Saudi Arabia, Kuwait, Oman, Qatar, and the UAE span 1990-2022, subject to availability. We estimate demand elasticities to assess the roles of price, industrial output, fuel substitution, and unobserved global/regional factors. Results show a positive income effect and a negative own-price effect across all countries, consistent with demand theory. Income elasticity ranges from 2.96 in Bahrain to 0.23 in Kuwait, with a panel value of 1.10. Own-price elasticity is inelastic, from 0.004 in Kuwait to 0.58 in Bahrain, with a panel value of 0.30. Cross-price effects vary: electricity prices raise gas demand in Algeria and Egypt, while gasoline prices matter in Algeria, Oman, and the UAE. Common global and regional factors significantly influence demand in Algeria, Saudi Arabia, Kuwait, and Oman.

C1436: Nonstationary panel approaches to approximating nonlinear steady state functions

Presenter: **Peter Pedroni**, Williams College, United States

The aim is to propose new panel time series methods for approximating a broad class of steady state functions that are of unknown form. In particular, the asymptotic and small sample properties of a number of potential estimation approaches which involve polynomial approximation methods are investigated. These include grouped, pooled, and time-averaged cross-sectional estimators. The data-generating processes are taken to be nonstationary and heterogeneous among units of the panel. Conditions required for the approximating function to converge to the mixture average of the true function are discussed, and small sample properties are studied via Monte Carlo simulations. An empirical illustration is provided for the environmental Kuznets curve.

CO414 Room BCB M201 INNOVATIVE STATISTICAL APPLICATIONS

Chair: Soudeep Deb

C0505: Conditional copula models using loss-based Bayesian additive regression trees

Presenter: **Tathagata Basu**, Newcastle University, United Kingdom

Co-authors: Fabrizio Leisen, Cristiano Villa, Kevin Wilson

The aim is to present a novel semi-parametric Bayesian approach for modeling conditional copulas to understand the dependence structure between two random variables when it is influenced by a different covariate. The use of Bayesian additive regression trees is proposed to model the conditional copulas. A loss-based prior is specified for the BART model suggested by a prior study, which is designed to reduce the loss in information and complexity for tree misspecification, giving a parsimonious model that avoids over-fitting, a common issue of BART models. Results are presented with both simulated and a real dataset to show the applicability and efficiency of the method.

C0614: Modeling purpose-driven tourism to India: A Bayesian approach with cross-category dependence

Presenter: **Amrutha Seshagiri**, Indian Institute of Management Bangalore (IIMB), India

Co-authors: Soudeep Deb

Tourism is a key contributor to economic development, playing a significant role in boosting GDP and generating employment. Understanding its determinants can provide valuable insights for policy planning and sectoral growth. Inbound tourism is analyzed in India from 62 countries, focusing on two primary factors: The relative economic strength of tourists' countries and their freedom levels, both in comparison to India. Tourist arrivals are categorized by purpose, which are business, leisure, visiting friends and family, medical, and other reasons, to explore how economic and freedom factors influence each category. As these categories are likely interdependent, a Bayesian modeling framework is employed with an error covariance structure to capture such dependencies. This approach helps account for unobserved factors influencing multiple travel purposes, improving the robustness of our estimates. Results show that relative economic strength and currency exchange rates have a significant positive effect on inbound tourism across categories. However, relative freedom status does not show a statistically significant influence. Strong interdependence is also found between medical visits and visits by the Indian diaspora, suggesting many non-resident Indians return for healthcare. These insights have important implications for tourism promotion strategies, particularly in targeting economically stronger countries and those with large Indian expatriate populations.

C0867: Scalable spatial skew-Gaussian process models

Presenter: **Kapil Gupta**, Indian Institute of Management, Kozhikode, India

Spatial data often exhibit skewness and heavy tails that violate the Gaussian assumptions underpinning traditional geostatistical methods. While transformation-based approaches such as trans-Gaussian kriging attempt to address this, they often suffer from bias and theoretical limitations. A Bayesian spatial modeling framework is proposed using a skewed Gaussian process model. The full conditional posteriors are derived, their normalization constants are analyzed, and closed-form expressions are provided. The model enables scalable inference through tractable likelihood decomposition and efficient MCMC sampling. The proposed methodology is validated through comprehensive simulation studies, demonstrating improved predictive performance and robustness over existing approaches.

C0916: Exploration, confirmation, and replication: A two team cross-screening

Presenter: **Samrat Roy**, Indian Institute of Management Ahmedabad, India

The long-term consequences of unwanted pregnancies carried to term on mothers have not been much explored. Data is used from the Wisconsin Longitudinal Study (WLS), and a novel method is proposed, namely two-team cross-screening, to study the possible effects of unwanted pregnancies carried to term on various aspects of mothers' later-life mental health, physical health, economic well-being, and life satisfaction. The method, unlike existing approaches to observational studies, enables the investigators to perform exploratory data analysis, confirmatory data analysis, and

replication in the same study. This is a valuable property when there is only a single data set available with unique strengths to perform exploratory, confirmatory, and replication analysis. In two-team cross-screening, the investigators split themselves into two teams, and the data is split as well according to a meaningful covariate. Each team then performs exploratory data analysis on its part of the data to design an analysis plan for the other part of the data. The complete freedom of the teams in designing the analysis has the potential to generate new unanticipated hypotheses in addition to a prefixed set of hypotheses. Moreover, only the hypotheses that looked promising in the data each team explored are forwarded for analysis (thus alleviating the multiple testing problem). These advantages are demonstrated in the study of the effects of unwanted pregnancies on mothers' later life outcomes

CO128 Room BCB M202 MODELING, FORECASTING, AND POLICY ASSESSMENT: MACRO, FINANCE Chair: Mohammad Jahan-Parvar

C0248: Can modern theories of structural change fit business cycles data?

Presenter: **Loris Rubini**, University of New Hampshire, United States

The purpose is to investigate the ability of workhorse structural change models in accounting for the business cycle properties of an economy. Three different preferences specifications are considered from past studies, paired with standard sectoral production functions with random total factor productivity (TFP) shocks. In each case, preference parameters are estimated using long-run structural change data, and common TFP processes calibrated on observed relative prices. Main results can be summarized by: i) all models display a volatility of aggregate variables substantially lower than the data, but they account for a large fraction of the volatility of consumption relative to GDP; ii) at the sectoral level, only CLM accounts for a substantial fraction of absolute and relative volatility; iii) all models do reasonably well in accounting for the cyclical of aggregate GDP components; and iv) only HRV can account for the cyclical of sectoral variables.

C0261: Trend-cycle decomposition and forecasting using Bayesian multivariate unobserved components

Presenter: **Mohammad Jahan-Parvar**, Federal Reserve Board of Governors, United States

Co-authors: Charles Knipp, Pawel Szerszen

A generalized multivariate unobserved components model is proposed to decompose macroeconomic data into trend and cyclical components. The series is then forecasted using Bayesian methods. It is documented that a fully Bayesian estimation, which accounts for state and parameter uncertainty, consistently dominates out-of-sample forecasts produced by alternative multivariate and univariate models. In addition, allowing for stochastic volatility components in variables improves forecasts. To address data limitations, cross-sectional information is exploited, the commonalities across variables are used, and both parameter and state uncertainty are accounted for. Finally, it is found that an optimally pooled univariate model outperforms individual univariate specifications and performs generally closer to the benchmark model.

C0196: A machine learning methodology for daily assessment of bank health, interconnectedness, and systemic risk

Presenter: **Celso Brunetti**, Bocconi University and Federal Reserve Board, United States

Co-authors: Shawn Mankad, Jeffrey Harris

A novel methodology is proposed to estimate the portfolio composition of banks as a function of daily stock returns. Building on a model where individual bank balance sheets connect through common holdings, a constrained semi-non-negative matrix factorization problem is derived and solved, where the rows (corresponding to banks) of one latent matrix factor (representing asset holdings) are subject to probability constraints. Although banks report assets at low frequencies, estimating factorization over a rolling window allows for the derivation of daily estimates of bank portfolios. Estimates of asset holdings are validated by showing they closely match balance-sheet data reported in quarterly regulatory filings.

C0216: Macroprudential policies and credit volatility

Presenter: **Giovanni Trovato**, University of Rome Tor Vergata, Italy

Co-authors: Alessio Farcomeni, Lorenzo Carbonari, Cosimo cosimopetracchi

The purpose is to present a model for data reduction and provide time-fixed indicators for macroprudential policies. Using a panel of 119 countries from 2000 to 2015, the effectiveness of macroprudential policies in reducing volatility in private credit is empirically assessed. Unobserved heterogeneity among countries is an important factor. An econometric model is employed that accounts for this heterogeneity and it is documented that the impact of macroprudential policies on financial stability varies, leading to either deterioration or improvement, depending on the macroeconomic conditions of the country in which they are implemented.

CO302 Room BCB 307 ADVANCES IN STATISTICAL MACHINE LEARNING

Chair: Tiffany Tang

C0350: Peer effects in the linear-in-means model may be inestimable even when identified

Presenter: **Alex Hayes**, Stanford University, United States

Co-authors: Keith Levin

Estimation in the linear-in-means model when a randomized treatment is applied to all nodes in a network, and the potential for an identifiability-estimability gap is shown. When treatment is assigned independently of the network structure, peer effects are identified but potentially inestimable due to an asymptotic collinearity issue. The estimation error is lower-bounded for ordinary least squares, and it is shown that these estimates may be inconsistent or fail to achieve nominal coverage rates whenever the harmonic mean degree of the network diverges with sample size. Simulations show that two-stage least squares and quasi-maximum likelihood estimators behave similarly. Results thus suggest caution when using the linear-in-means model to model spillovers in random experiments on dense networks. The behavior of the linear-in-means model is further investigated when covariates are endogenous and associated with network structure. It is shown that explicitly modeling homophily with random dot product graphs can prevent asymptotic collinearity and estimability issues, provided that there is sufficient degree heterogeneity in the network.

C0556: Factor adjusted spectral clustering for mixture models

Presenter: **Soham Jana**, University of Notre Dame, United States

Co-authors: Shange Tang, Jianqing Fan

A factor modeling-based approach is studied for clustering high-dimensional data generated from a mixture of strongly correlated variables. Standard techniques for clustering high-dimensional data, e.g., naive spectral clustering, often fail to yield insightful results as their performances heavily depend on the mixture components having a weakly correlated structure. To address the clustering problem in the presence of a latent factor model, the factor-adjusted spectral clustering (FASC) algorithm is proposed, which uses an additional data denoising step via eliminating the factor component to cope with the data dependency. This method is proven to achieve an exponentially low mislabeling rate, with respect to the signal-to-noise ratio, under a general set of assumptions. The assumption bridges many classical factor models in the literature, such as the pervasive factor model, the weak factor model, and the sparse factor model. The FASC algorithm is also computationally efficient, requiring only near-linear sample complexity with respect to the data dimension. The applicability of the FASC algorithm is also shown with real data experiments and numerical studies, and it is established that FASC provides significant results in many cases where traditional spectral clustering fails.

C0921: Valid f-screening in linear regression

Presenter: **Daniel Kessler**, University of North Carolina at Chapel Hill, United States

Co-authors: Olivia McGough, Daniela Witten

Suppose that a data analyst wishes to report the results of a linear regression only if the overall null hypothesis is rejected. This practice, which is referred to as F-screening, is in fact common practice across a number of applied fields. Unfortunately, it poses a problem: Standard guarantees

for the inferential outputs of linear regression, such as Type 1 error control of hypothesis tests and nominal coverage of confidence intervals, hold unconditionally, but fail to hold conditional on rejection of the overall null hypothesis. An inferential toolbox is developed for the coefficients in a least squares model that are valid conditional on rejection of the overall null hypothesis. Selective p-values that lead to tests are developed that control the selective Type 1 error, i.e., the Type 1 error conditional on having rejected the overall null hypothesis. Furthermore, they can be computed without access to the raw data, i.e., using only the standard outputs of a least squares linear regression, and therefore are suitable for use in a retrospective analysis of a published study. Confidence intervals are also developed that attain nominal selective coverage, and point estimates that account for having rejected the overall null hypothesis. It is shown empirically that the selective procedure is preferable to an alternative approach that relies on sample splitting, and its performance is demonstrated via re-analysis of two datasets from the biomedical literature.

C0702: Patchwork PCA: Joint dimension reduction for semi-overlapping data patches

Presenter: **Lili Zheng**, University of Illinois Urbana - Champaign, United States

Patchwork learning arises as a new and challenging data collection paradigm where both samples and features are observed in fragmented subsets. Due to technological limits, measurement expense, or multimodal data integration, such patchwork data structures are frequently seen in neuroscience, healthcare, and genomics, among others. Instead of analyzing each data patch separately, it is highly desirable to extract comprehensive knowledge from the whole data set. A new PCA method designed for patchwork learning is introduced, which extracts principal components for the whole sample and feature space based on a collection of semi-overlapping data patches. It is demonstrated how key challenges are addressed for patchwork data, such as non-random missingness, heterogeneous SNRs, irregular observational patterns, etc. Statistical error bounds are shown for the estimated principal components and sample loadings, as well as demonstrating their performance on real biomedical data sets.

CO272 Room BCB 308 APPLICATION OF MACHINE LEARNING IN SAMPLE SURVEYS AND SMALL AREA ESTIMATION Chair: Aditi Sen

C0475: Gradient boosting for hierarchical data in small area estimation

Presenter: **Paul Messer**, University of Bamberg, Germany

Co-authors: Timo Schmid

Small area estimation (SAE) combines survey data with auxiliary sources such as administrative records, census data, or alternative datasets that typically offer broader coverage. By integrating these sources, SAE enhances the accuracy of (direct) survey estimates. To account for the hierarchical structure of survey data, model-based SAE methods often rely on linear mixed models (LMMs). However, the distributional (e.g., normality) and structural (e.g., linearity) assumptions of LMMs may not always hold in practice, and the accuracy of model-based SAE depends on the validity of these assumptions. To address these limitations, a mixed-effect gradient boosting (MEGB) approach is proposed, which combines the flexibility of gradient boosting machines with the ability of mixed models to account for hierarchical data structures. MEGB extends standard gradient boosting by incorporating random effects, allowing it to capture unobserved heterogeneity across domains while retaining a nonparametric framework that models non-linearities and interactions in the data. MEGB supports the derivation of area-level means from unit-level data and uses a nonparametric bootstrap to estimate the mean squared error. Its performance is assessed through a model-based simulation study, comparing MEGB to established estimators, and further demonstrated using real-world data. The results suggest that MEGB offers promising area mean estimates and may outperform existing SAE methods in various scenarios.

C0529: Evaluating alternative deep learning approaches for village-level wealth estimation using satellite imagery

Presenter: **David Newhouse**, World Bank Group, United States

Co-authors: Diana Jaganjac, Josh Merfeld, Kushan Weerakoon

Building on existing research that uses satellite imagery and auxiliary data to estimate poverty at hyperlocal levels, transformer architectures and convolutional neural networks are evaluated to generate estimates of mean asset index values at the enumeration area level in Malawi. Estimates are generated using Planet Imagery and evaluated in held-out test sets after combining two household surveys: The Integrated Household Survey and the Multiple Indicator Cluster Survey. Estimates generated using Resnet 18 and the first version of the Prithvi foundational model outperform other architectures, achieving out of sample Pearson correlations of approximately 0.81. This exceeds performance from the recently developed ConvNeXt convolutional neural network and two standard vision transformer models. A robustness check using the Prithvi model with lower-resolution Landsat imagery achieves an out-of-sample correlation of 0.71. The results indicate that large-scale utilization of new foundational models to combine household survey data and satellite imagery offers a promising approach to generating accurate village-level estimates of wealth indices in this context.

C0589: Multi-target semi-supervised learning with application to small area estimation

Presenter: **Katarzyna Reluga**, Humboldt University of Berlin, Germany

Co-authors: Nicola Salvati, Mark van der Laan

In the classical single-target semi-supervised learning (SSL) setting, one has access to (i) a moderately sized labeled dataset containing both response values and associated features, and (ii) a much larger unlabeled dataset with only covariates observed. SSL naturally arises in settings where collecting features is easy, but obtaining labels is expensive or time-consuming, for example, in electronic health records or survey data, where full data is available for only a small subset of the population. This framework is extended to multi-target semi-supervised learning, where the goal is to estimate several parameters of interest across different subpopulations, but labeled data are sparse. Classical SSL methods can suffer from excessive variability in this setting. Novel estimation methods tailored to this problem are proposed, and it is demonstrated how they improve stability and efficiency. Finally, it is shown how small area estimation emerges as a special case of this broader learning framework.

C0775: Improving measurement error and representativeness in nonprobability surveys

Presenter: **Aditi Sen**, University of Maryland, College Park, United States

Co-authors: Partha Lahiri

In the age of big data, nonprobability surveys are becoming increasingly abundant. Data integration techniques involving both probability and nonprobability surveys are being extensively used for providing improved estimates for finite population estimation. While much of the existing research has focused on mitigating selection bias in nonprobability surveys, the issue of measurement error within these surveys remains relatively unexplored. Statistical methods devised with the purpose of reducing selection bias are appropriate for reliable estimation, only under the assumption of the accuracy of survey responses. Motivated by a recent case study, the research addresses bias from both measurement and sampling errors in nonprobability surveys. A new data integration method is proposed that leverages machine learning models to construct a composite estimator. The composite estimator integrates probability and nonprobability surveys when both contain response variables of interest. The performance of this estimator is analyzed in comparison to an existing composite estimator in literature, analytically as well as empirically, using multiple survey data from a recent study. Finally, conditions are identified under which the proposed estimator outperforms estimators based solely on probability surveys.

CO030 Room BCB 309 ADVANCES IN HIGH-DIMENSIONAL TIME SERIES AND BAYESIAN MODELING Chair: S Yaser Samadi

C1429: A semi-parametric approach for clustering high-dimensional, non-stationary, auto-correlated time series

Presenter: **Qiyuan Wang**, Texas A&M University, United States

A novel semi-parametric estimation algorithm is developed for accurately estimating time-varying mean and variance in autoregressive (AR) models. Utilizing B-splines with generalized least squares (GLS) estimation for smooth parametrization and weighted least squares (WLS) for more

precise estimation, the approach addresses the challenges posed by time-varying dynamics in time series data. The covariance matrix in the GLS estimation of the spline coefficients is iteratively updated by calculating it through the WLS estimation of the AR coefficients in a band-limited manner. Meanwhile, a new autoregressive model is proposed that incorporates time-varying variance with a finite bounded envelope function, and a novel method is introduced to estimate it through splines. Additionally, the order of the AR model is determined through a generalized Bayesian information criterion (GBICp) that incorporates prior information. The effectiveness of the methodology is demonstrated through extensive simulations and applications to real-world electrocardiograms (ECGs) data, showcasing significant improvements in the dimension reduction while preserving major features for high-accuracy clustering tasks.

C1445: Bayesian copula factor autoregressive models for time series mixed data

Presenter: **S Yaser Samadi**, Southern Illinois University Carbondale, United States

Co-authors: Samira Zaroudi, Hadi Safari-Katesari

The aim is to propose a Bayesian copula factor autoregressive (BCFAR) model for analyzing time series mixed data, accommodating main effects and interactions. The main motivation is to infer dynamic interactions between macroeconomic variables and stock market indices. The BCFAR model assumes conditional independence and applies latent factors in both response time series and high-dimensional mixed-type covariates in quadratic regression using copula functions. To complement this, a simpler time series Bayesian factor regression (TS-BFR) model is introduced, tailored for continuous Gaussian multivariate time series. Both models build on the quadratic autoregression (QAR) framework, employ latent factors for efficient dimension reduction, and capture main effects and interactions of covariates by integrating latent variables into the response. For computational efficiency, a semiparametric time series extended rank likelihood is used for explanatory-variable margins in the BCFAR model, reducing parameters and ensuring fast computation. To estimate latent factors and parameters, flexible Bayesian algorithms are designed, employing Metropolis-Hastings (MH) and forward filtering backward sampling (FFBS) within Gibbs sampling. The effectiveness of these methods is shown through simulation studies, and the approach is further validated with a macroeconomic dataset.

C1446: Fast, efficient, and automatic tuning parameter selection for LASSO

Presenter: **Sumanta Basu**, Cornell University, United States

Tuning parameter selection for penalized regression methods such as LASSO is an important issue in practice, albeit less explored in the literature of statistical methodology. Most common choices include cross-validation (CV), which is computationally expensive, or information criteria such as AIC or BIC, which are known to perform worse in high-dimensional scenarios. Guided by the asymptotic theory of LASSO that connects the choice of tuning parameter to estimation of error standard deviation, autotune is proposed, a procedure that alternately maximizes a (restricted) penalized log-likelihood over regression coefficients and the nuisance parameter, resulting in an automatic tuning algorithm. The core insight behind autotune is that under exact or approximate sparsity conditions, estimation of the scalar nuisance parameter may often be statistically and computationally easier than estimation of the high-dimensional regression parameter, leading to a gain in efficiency. Using simulated and real data sets, it is shown that autotune is faster, and provides superior estimation, variable selection, and prediction performance than existing tuning strategies for LASSO as well as alternatives such as the scaled LASSO. The algorithm can be extended naturally to high-dimensional time series problems, and this is illustrated in the context of estimating large vector autoregression (VAR).

C1450: Dimension reduction in VAR models via informative lag selection

Presenter: **Wiranthe Herath**, Drake University, United States

Co-authors: S Yaser Samadi

The increasing dimensionality of multivariate time series data creates significant challenges for traditional vector autoregressive (VAR) models, frequently resulting in overfitting, inefficient estimation, and poor forecasting ability. To address these issues, two novel VAR models are developed that allow for targeted dimension reduction by focusing on the most informative lagged predictors. These models introduce a simple structure that not only improves estimation efficiency but also forecasting accuracy compared to traditional VAR approaches. Extensive simulation results and empirical analyses using finance and macroeconomic data demonstrate consistent performance gains in both parameter estimation and predictive outcomes, emphasizing the practical value of our techniques in multivariate time series contexts.

CO326 Room BCB 310 STATISTICAL METHODOLOGY FOR DATA WITH SPATIAL AND TEMPORAL DEPENDENCIES Chair: Hyebin Song

C0262: Topological signal processing

Presenter: **Robin Belton**, Vassar College, United States

The aim is to study a popular tool in topological data analysis (TDA) called sublevel set persistent homology on discrete functions through the perspective of finite ordered sets of both linearly ordered and cyclically ordered domains. The duality of filtrations of sublevel sets is proven, which undergirds a range of duality results of sublevel set persistent homology without the need to invoke continuous functions or classical Morse theory. Furthermore, aspects of the surgery of circular and linearly ordered sets are discussed, with a focus on applications in audio signal processing. It ends by discussing ideas of future work that integrate this framework with more traditional statistical techniques for analyzing time series.

C0572: Bandwidth selection for zero Lugsail estimators

Presenter: **Rebecca Kurtz-Garcia**, Smith College, United States

Test statistics, confidence intervals, and p-values all typically rely on an estimate for variance. For data sets that are not independent and identically distributed (iid), caution must be used when selecting a variance estimator. If the dependence structure is unknown but stationary, a robust long-run variance (LRV) estimator can be used, which can handle a wide variety of scenarios. Spectral variance (SV) estimators are one of the most common LRV estimation methods, but they suffer from a negative bias in the presence of positive correlation. An alternative zero lugsail estimator has been proposed to combat this issue, which has a zero asymptotic bias regardless of correlation. Both SV and zero lugsail estimators rely on a bandwidth parameter, a critical component for the estimation process. Currently, no guidelines exist for selecting a bandwidth for the zero lugsail estimator. An optimal bandwidth rule is proposed for zero lugsail estimators when relying on nonstandard limiting distributions. With this procedure, bias can be greatly improved, variability accounted for, and an estimator optimized for inference obtained.

C0610: Estimating the effect of spatial interventions in the presence of interference: A causal inference approach

Presenter: **Nathan Wikle**, University of Iowa, United States

In many settings, it is of primary interest to estimate the effect of a spatially varying intervention that has a nonlocal (i.e., spillover) effect on its surrounding environment. For example, what is the effect of a coal-fired power plant on downwind air pollution concentrations, or how might the addition of a rural health clinic improve health outcomes for individuals living within a certain distance? Unfortunately, estimating causal effects from observational data in such settings is challenging, due to (i) the risk of confounding bias, and (ii) the potential for treatment interference, namely, that multiple interventions affect the same outcome locations. A framework is introduced for causal inference with spatial data in which causal estimands are defined as functionals of the potential outcome distribution under a set of stochastic interventions. Corresponding nonparametric identifying assumptions are considered, allowing the estimands to be estimated from observational data in the presence of distance-limited interference, and an augmented inverse probability of treatment-type estimator is proposed. Notably, the estimator is constructed from a log-Gaussian Cox process model for intervention locations and a semiparametric outcome model of the spillover structure that accounts for spatial autocorrelation. The proposed method is used to estimate the effect of large concentrated animal feeding operations (CAFOs) on the environment in Iowa.

C1054: Non-parametric mixture models for covariance function estimation**Presenter:** Stephen Berg, Penn State University, United States

An approach for estimating covariance functions based on nonparametric mixture models will be introduced, with an emphasis on a weighted least squares estimator of the autocovariance sequence from a reversible Markov chain. The estimator is shown to lead to strongly consistent estimates of the asymptotic variance of the sample mean from an MCMC sample, as well as to consistent estimates of the autocovariance sequence. An algorithm for computing the estimator is presented, and some empirical applications will be shown. The proposed shape-constrained estimator exploits a mixture representation of the autocovariance sequence from a reversible Markov chain. Similar mixture representations exist for stationary covariance functions in spatial statistics, including for the Matern covariance as a special case, and some extensions of shape-constrained approaches are highlighted for estimating covariance functions in spatial statistics.

CO113 Room BCB 311 ROBUSTNESS AND REGULARIZATION IN MIXTURE MODELING AND NETWORK ANALYSIS Chair: Suyeon Kang**C0306: Robust estimation and outlier detection for stochastic block models via subgraph search****Presenter:** Christine Keribin, INRIA-Paris-Saclay University, France**Co-authors:** Leonardo Martins Bianco, Zacharie Naulet

Community detection is a fundamental task in graph analysis, with methods often relying on fitting models like the stochastic block model (SBM) to observed networks. While many algorithms can accurately estimate SBM parameters when the input graph is a perfect sample from the model, real-world graphs rarely conform to such idealized assumptions. Therefore, robust algorithms are crucial ones that can recover model parameters even when the data deviates from the assumed distribution. SubSearch is proposed, an algorithm for robustly estimating SBM parameters by exploring the space of subgraphs in search of one that closely aligns with the model's assumptions. The approach also functions as an outlier detection method, properly identifying nodes responsible for the graph's deviation from the model and going beyond simple techniques like pruning high-degree nodes. Extensive experiments on both synthetic and real-world datasets demonstrate the effectiveness of the method.

C0825: Robust finite mixture of regression model selection**Presenter:** Frans Kanfer, University of Pretoria, South Africa**Co-authors:** Andre Kleynhans, Sollie Millard

Finite mixture of regression (FMR) models are widely employed for analysing data observed from heterogeneous subpopulations. Despite the flexibility of FMR models, estimation procedures face challenges, such as pre-specification of the number of components, order selection, and the identification of informative variables and covariates. Standard estimation procedures are also often sensitive to non-typical observations, affecting estimated model performance. An approach is proposed that integrates a penalized information criterion with a self-paced learning (SPL) algorithm. The penalization mechanism enables joint order and variable selection, while the SPL framework regulates the inclusion of observations during the learning process, thereby mitigating the influence of non-typical data. The proposed method is assessed through simulation studies, evaluating its robustness in component number estimation, variable selection, and parameter recovery under varying levels of contamination of non-typical observations. Real-world applications are also considered.

C1406: A Bayesian nonparametric approach to discriminant analysis**Presenter:** Bernardo Nipoti, University of Milan Bicocca, Italy**Co-authors:** Laura D Angelo, Tommaso Rigon

A Bayesian nonparametric framework is introduced to improve classical discriminant analysis, particularly in scenarios with sparse data. The method provides a flexible approach that encompasses both linear and quadratic discriminant analysis as special cases. The key innovation lies in allowing information sharing across groups to improve the estimation of group-specific covariance matrices. This is accomplished through a scale-only nonparametric mixture model defined on the space of positive definite matrices. The use of a conjugate nonparametric prior ensures tractability and remarkable ease of implementation. Applications to both simulated and real datasets demonstrate the adaptability and effectiveness of the proposed methodology.

C0957: Reduced-rank finite mixture regression for multivariate response via low-rank regularization**Presenter:** Suyeon Kang, University of Central Florida, United States**Co-authors:** Kun Chen, Weixin Yao

Given the rapid growth in data volume and access to diverse data sources, data complexity and heterogeneity have escalated across many fields. The aim is to extend reduced-rank estimation to mixture modeling by proposing a new class of reduced-rank multivariate mixture regression models. These models handle multiple continuous responses under population heterogeneity while extracting low-dimensional structure. Computationally efficient EM-type algorithms that incorporate both a rank penalty and an adaptive nuclear-norm penalty are derived, enabling simultaneous subgroup identification, parameter estimation, and rank selection. The monotonicity of the penalized likelihood sequence and the asymptotic consistency of the estimators are proven. Simulation studies and real data analysis have been carried out to validate the effectiveness and practical usefulness of the proposed methods. The R package rrMixture is developed for the implementation and is publicly available on CRAN.

CO336 Room BCB 312 STATISTICAL GENETICS AND MULTIOMIC DATA ANALYSIS Chair: Indranil Mukhopadhyay**C0313: Scalable Bayesian multivariate regression analysis for selecting targeted regressors in microbiome analysis****Presenter:** Sounak Chakraborty, University of Missouri, Columbia, United States**Co-authors:** Priyam Das, Tanujit Dey, Christine Peterson

B-MASTER (Bayesian multivariate regression analysis for selecting targeted essential regressors) is introduced, a fully Bayesian framework for scalable multivariate regression in high dimensions. B-MASTER is designed to identify master predictors covariates exerting widespread influence across many outcomes via a hybrid penalty: An ℓ_1 penalty induces elementwise sparsity, while an ℓ_2 penalty enforces groupwise shrinkage across rows of the coefficient matrix. This structure selects a parsimonious set of key covariates, enhancing interpretability. A tailored Gibbs sampler achieves scalability, with runtime growing linearly in parameter dimension and remaining stable across sample sizes; full posterior inference is feasible for models with up to four million parameters. Posterior consistency and contraction rate results are established, showing that B-MASTER concentrates around the truth at the minimax-optimal rate under sparsity. These theoretical guarantees are supported by strong empirical performance: In simulations, B-MASTER outperforms competing methods in estimation and signal recovery. Applied to microbiome metabolomics data from colorectal cancer patients, B-MASTER reveals microbial genera that shape broad metabolite profiles, uncovering relationships missed by other methods. The proposed approach is principled, interpretable, and scalable for discovering systemic patterns in ultra-high-dimensional biomedical data.

C0736: Integrative multi-omics QTL colocalization maps regulatory architecture in aging human brain**Presenter:** Kushal Dey, Memorial Sloan Kettering Cancer Center, United States

Multi-trait QTL (xQTL) colocalization has shown great promise in identifying causal variants with shared genetic etiology across multiple molecular modalities, contexts, and complex diseases. ColocBoost is proposed, a multi-task learning colocalization method that can scale to hundreds of traits while accounting for multiple causal variants. ColocBoost employs a specialized gradient boosting framework that can adaptively couple colocalized traits while performing causal variant selection, thereby enhancing the detection of weaker shared signals compared to existing pairwise and

multi-trait colocalization methods. ColocBoost is applied genome-wide to 17 gene-level single-nucleus and bulk xQTL data from the aging brain cortex of ROSMAP individuals (average $N = 595$), encompassing 6 cell types, 3 brain regions and 3 molecular modalities (expression, splicing, and protein abundance). Across molecular xQTLs, ColocBoost identified 16,503 distinct colocalization events, exhibiting 10.7(0.74)-fold enrichment for heritability across 57 complex diseases/traits and showing strong concordance with element-gene pairs validated by CRISPR screening assays. When colocalized against Alzheimer's disease (AD) GWAS, ColocBoost identified up to 2.5-fold more distinct colocalized loci. Several genes, like BLNK and CTSB, showed sub-threshold associations in GWAS but were identified through multi-omics colocalizations showing functional involvement in AD pathogenesis.

C0876: **Robust high-dimensional variable selection for GWAS via the Bayesian group lasso**

Presenter: **Mayetri Gupta**, University of Glasgow, United Kingdom

Co-authors: lanxin li, Vincent Macaulay, Indranil Mukhopadhyay

Genome-wide association studies (GWAS) are a powerful tool for exploring the connections between human disease and genetic variation. The accuracy and reproducibility of association detection in such studies are typically limited due to the weakness of association signals and local correlations between genetic variants. Set-based association tests often tend to be used for their ability to aggregate weak signals and alleviate multicollinearity issues. Statistical methods proposed for set-based association analyses, however, mainly focus on testing the marginal association of each set with a trait at a time, ignoring information across the genome. To improve detection accuracy, the proposal is to examine a set of groups of correlated genetic variants simultaneously via a Bayesian modelling framework, adapting ideas from Bayesian group lasso for high-dimensional regression. In the proposed model, a double-shrinkage hierarchical prior captures the sparsity pattern in GWAS, while a population-based Markov-chain Monte Carlo sampler is used for efficient posterior sampling. The implementation is further sped up by a group-based split-and-merge model-fitting strategy. The model appears to be powerful in locating true association signals both at a group and individual variant level in simulations and a GWAS from an Alzheimer's disease study, and offers improved performance in terms of sensitivity and precision compared to other existing methods.

C1131: **Bayesian meta-analysis of penetrance for cancer risk with an extension to include studies with ascertainment bias**

Presenter: **Swati Biswas**, University of Texas at Dallas, United States

Co-authors: Thanthirige Lakshika M Ruberu, Danielle Braun, Giovanni Parmigiani

Multi-gene panel testing allows many cancer susceptibility genes to be tested quickly at a lower cost, making such testing accessible to a broader population. Thus, more patients carrying pathogenic germline mutations in various cancer-susceptibility genes are being identified. This creates both opportunity and urgency to provide appropriate risk-reducing guidance, which depends on accurate age-specific cancer risk (penetrance) estimates for each gene. A meta-analysis approach is proposed using a Bayesian hierarchical random-effects model to estimate penetrance by integrating studies that report various risk measures (e.g., penetrance, relative risk, odds ratio), while accounting for uncertainty. Using a Markov chain Monte Carlo algorithm, posterior distributions are derived to estimate penetrance and credible intervals. The method is assessed through simulations involving risk estimates for two moderate-risk breast cancer genes, ATM and PALB2, and superior performance is demonstrated over existing methods in terms of coverage probability and mean squared error. The model is also extended to account for ascertainment bias by incorporating a bias term with appropriate priors. Simulations show that adjusting for this bias leads to more accurate and precise penetrance estimates than ignoring it or discarding biased studies. Finally, the method is applied to estimate breast cancer penetrance in carriers of pathogenic variants in ATM and PALB2.

CO325 Room BCB 313 ADVANCES IN SEQUENTIAL DECISION MAKING

Chair: Sakshi Arya

C0186: **Smooth contextual bandits: Bridging the parametric and non-differentiable regret regimes**

Presenter: **Yichun Hu**, Cornell University, Johnson Graduate School of Management, United States

A nonparametric contextual bandit problem is studied in which the expected reward functions belong to a Holder class with smoothness parameter β . It is shown how this interpolates between two extremes that were previously studied in isolation: Nondifferentiable bandits (β at most 1), with which rate-optimal regret is achieved by running separate noncontextual bandits in different context regions, and parametric-response bandits (infinite β), with which rate-optimal regret can be achieved with minimal or no exploration because of infinite extrapolatability. A novel algorithm is developed that carefully adjusts to all smoothness settings, and its regret is proven to be rate-optimal by establishing matching upper and lower bounds, recovering the existing results at the two extremes. In this sense, the gap is bridged between the existing literature on parametric and nondifferentiable contextual bandit problems and between bandit algorithms that exclusively use global or local information, shedding light on the crucial interplay of complexity and regret in contextual bandits.

C0483: **On stopping times of power-one sequential tests: Tight lower and upper bounds**

Presenter: **Shubhada Agrawal**, Indian Institute of Science Bangalore, India

Co-authors: Aaditya Ramdas

Two lower bounds are proven for stopping times of sequential tests between general composite nulls and alternatives. The first lower bound is for the setting where the type-1 error level α approaches zero, and equals $\log(1/\alpha)$ divided by a certain infimum KL divergence, termed KL-inf. The second lower bound applies to the setting where α is fixed and KL-inf approaches 0 (meaning that the null and alternative sets are not separated) and equals $c \text{KL-inf} \log \log(\text{KL-inf})$ for a universal constant $c > 0$. A sufficient condition is also provided for matching the upper bounds, and it is shown that this condition is met in several special cases. Given past work, these upper and lower bounds are unsurprising in their form; the main contribution is the generality in which they hold, for example, not requiring reference measures or compactness of the classes.

C0488: **Semiparametric contextual bandits with sufficient dimension reduction**

Presenter: **Sakshi Arya**, Case Western Reserve University, United States

A novel semi-parametric framework is introduced for batched contextual multi-armed bandits that leverages a single-index regression model to flexibly capture relationships between covariates and arm rewards. The proposed algorithm, batched single-index dynamic binning and successive arm elimination (BIDS), combines dynamic binning based on the estimated single-index direction with a successive arm elimination strategy. This approach accommodates both settings where a pilot index direction is known and where it must be estimated from data. For both cases, theoretical regret guarantees are derived, and it is shown that, when the single-index direction is estimated with sufficient accuracy, BIDS achieves minimax-optimal regret rates comparable to nonparametric bandits with a one-dimensional covariate, thereby circumventing the curse of dimensionality. Extensive experiments on simulated and real-world datasets demonstrate that BIDS outperforms existing nonparametric batched bandit methods in both sample efficiency and empirical performance, establishing the practical value of leveraging single-index structures in batch decision-making.

C0654: **Structured expert judgment for sequential decision-making**

Presenter: **Tina Nane**, TU Delft, Netherlands

Expert opinion is typically employed to support decision-making when data are not (or no longer) representative or simply unavailable, in high uncertainty contexts. The classical model for structured expert judgment (SEJ) provides a structured and validated methodology to elicit and mathematically aggregate experts' uncertainty assessments. Probability distributions are constructed from the (typically three) elicited percentiles. Uncertainty assessments of calibration variables, for which realizations are known to the analyst but not to the experts, enable a validation step, where objective measures indicate how statistically accurate and informative the assessments are. These measures yield performance-based weights, which are used to aggregate experts' assessments into distributions for variables of interest. The classical model for SEJ has been applied in more

than 250 professional studies across numerous domains, spanning from epidemiology, natural hazards, and climate change, to energy (transition). Most of the studies are a one-time application of SEJ. The aim of this research is to extend the classical model for SEJ to a sequential context. Through the dynamic nature of the data, when realizations become quickly available, the framework enables the sequential validation of individual and aggregated assessments. Performance metrics are derived to support sequential decision-making. Results from two studies are presented.

C0933: Adaptive partial monitoring in non-stationary environments

Presenter: **Henry Reeve**, Nanjing University, China

A framework is introduced for sequential decision making in a non-stationary stochastic environment in which the outcome distribution may change over time and the rewards may not be directly observable. The regret of a policy contrasts performance with the expected reward of a dynamic oracle capable of selecting an optimal sequence of actions for the non-stationary stochastic environment. An algorithm is provided which leverages e-processes to provably adapt to distributional changes in settings where the reward attained from a given action is not directly observed. It is found that the optimal regret depends upon a fascinating interplay between the level of observability, the noise level, the complexity of the action space, and the degree of non-stationarity.

CO158 Room BCB 402 NOVEL STATISTICAL AND COMPUTATIONAL METHODS FOR BIOMEDICAL SCIENCES

Chair: Fan Bu

C0204: Regularizing extrapolation in causal inference

Presenter: **Harsh Parikh**, Yale University, United States

Co-authors: Avi Feller, Elizabeth Stuart, Kara Rudolph, David Arbour

Linear smoothers in machine learning and causal inference predict using weighted averages of training outcomes. Traditional approaches either allow negative weights (improving feature balance but increasing variance and model dependence) or restrict weights to be non-negative (reducing variance but worsening imbalance). Replacing hard non-negativity constraints are proposed with soft penalties on extrapolation, introducing a "bias-bias-variance" tradeoff balancing feature imbalance, model misspecification, and estimator variance. A worst-case extrapolation error bound is derived, and an optimization procedure that regularizes extrapolation while minimizing imbalance is developed. The framework unifies existing methods, allowing practitioners to navigate the tradeoff with a single hyperparameter controlling the extrapolation penalty. Synthetic experiments demonstrate that the approach achieves better bias-variance tradeoffs than existing methods across various degrees of positivity violation and model misspecification. In a real-world application, generalizing randomized trial results to target populations, the method produces more reliable estimates while quantifying sensitivity to modeling assumptions. This enables researchers to make informed decisions about appropriate extrapolation levels rather than being restricted to either unconstrained or strictly non-negative weighting.

C0374: ProjMCMC: Scalable and stable posterior inference for Bayesian spatial factor models

Presenter: **Lu Zhang**, University of Southern California, United States

Factor models exhibit a fundamental tradeoff among flexibility, identifiability, and computational efficiency. Bayesian spatial factor models, in particular, face pronounced identifiability concerns and scaling difficulties. To mitigate these issues and enhance posterior inference reliability, projected Markov chain Monte Carlo (ProjMC²) is proposed, a novel Markov chain Monte Carlo (MCMC) sampling algorithm employing projection techniques and conditional conjugacy. ProjMC² is showcased within the context of spatial factor analysis, significantly improving posterior stability and MCMC mixing efficiency by projecting posterior sampling of latent factors onto a subspace of a scaled Stiefel manifold. Theoretical results establish convergence to the stationary distribution irrespective of initial values. Integrating this approach with scalable univariate spatial modeling strategies yields a stable, efficient, and flexible modeling and sampling methodology for large-scale spatial factor models. Simulation studies demonstrate the effectiveness and practical advantages of the proposed methods. The practical utility of the methodology is further illustrated through an analysis of spatial transcriptomic data obtained from human kidney tissues, showcasing its potential for enhancing the interpretability and robustness of spatial transcriptomics analyses.

C0375: Dependent Dirichlet-multinomial processes with random number of components

Presenter: **Andrea Cremaschi**, IE University, Spain

Co-authors: Beatrice Franzolini

Bayesian nonparametric methods under partial exchangeability have largely focused on infinite support priors, yet almost-surely finite dimensional dependent mixtures remain underexplored despite their strong theoretical guarantees and performance in the exchangeable case. A novel class of finite dependent Dirichlet-multinomial processes and their counterparts are introduced, incorporating a prior on the number of components. The class is built on generalized Wishart unnormalized weights. It is proven that the normalized diagonals of correlated Wishart matrices admit a hierarchical negative binomial Dirichlet representation whose marginal laws are Dirichlet. The key distributional properties of the induced random probability measures are derived, showing they can achieve any prescribed dependence level and support efficient posterior computation without the need for costly variable augmentation schemes. The practical advantages of the proposed processes are further demonstrated through extensive simulation studies and the analysis of two real datasets: A benchmark dataset and a case study on sex-specific gene expression differences in the human brain, highlighting its flexibility and computational efficiency in capturing complex dependence structures.

C0582: Hierarchical Bayesian framework for evidence synthesis across federated data networks

Presenter: **Fan Bu**, University of Michigan, United States

Co-authors: Weixiong Hua

Evidence synthesis, or meta-analysis, is crucial for combining results from multi-site studies, especially when subject-level data cannot be shared. Traditional meta-analysis approaches relying on mixed effects models with normality assumptions are often inadequate when single sites have small sample sizes or rare outcome events. In addition, bias can arise due to systematic error when using observational data sources. To address both challenges, a Bayesian evidence synthesis framework is proposed that performs likelihood-based meta-analysis and meanwhile corrects for bias. A Bayesian hierarchical model that jointly analyzes profile likelihoods instead of point estimates is introduced, which well accommodates small sample sizes and rare outcomes. Bias correction is achieved by learning an empirical bias distribution from a large number of negative control outcomes that are not associated with exposures of interest. A hierarchical Dirichlet process prior is employed to flexibly characterize bias globally and locally while borrowing information across data sources. A Hamiltonian Monte Carlo algorithm is implemented to efficiently carry out posterior inference. Through a series of simulation studies and a case study on comparing Type-2 diabetes treatments across a multi-national data network, it is demonstrated that the approach can provide more precise estimates and coherent results for improved interpretability.

CO281 Room BCB 403 ADVANCED STATISTICAL METHODS FOR SPATIAL TRANSCRIPTOMICS DATA ANALYSIS

Chair: Thierry Chekouo

C0533: Network models for spatial transcriptomics data

Presenter: **Satwik Acharyya**, University of Alabama at Birmingham, United States

Network models are powerful tools to investigate complex dependence structures in high-throughput genomic datasets. They allow for a holistic, systems-level view of the various biological processes, for intuitive understanding and coherent interpretations. However, most existing network or graphical models are developed under assumptions of homogeneity of samples and are not readily amenable to modeling spatial heterogeneity, which often manifests in spatial genomics data. Two spatial network models are discussed, focusing on spatially varying covariance and precision matrices. (I) SpaceX (spatially dependent gene co-expression network) is a Bayesian methodology to identify both shared and cluster-specific

co-expression networks across genes. (II) Spatial Graphical Regression (SGR) is a flexible approach based on graphical regression that enables spatially varying graphs over the spatial domain of the tissue. The framework incorporates multiple spatial covariates and provides a linear and non-linear functional mapping between the spatial domain and the precision matrices. All the approaches are illustrated by using case studies from cancer genomics.

C0747: Spatial omics driven computational discovery of prognostic cellular neighborhoods in tumor ecosystems

Presenter: Lulu Shang, MD Anderson, United States

Recent advancements in spatial omics have deepened the understanding of tissue microenvironments and cellular neighborhoods in tumors and other diseases. Existing neighborhood detection methods lack a universal, interpretable framework, yielding ambiguous neighborhoods that cannot be aligned across patients, disease stages, or platforms, hindering clinical translation. The aim is to present SCOPE (Systematic prOfiling of disease-relevant cellular Ecotypes), which extracts functional spatial ecotypes directly from single-cell spatial data and scales to thousands of tissue sections and millions of cells. SCOPE is benchmarked on both simulated and real datasets, showing that it can identify ecotypes with known prognostic patterns more accurately than existing methods. Ecotype enrichment is found across patient subtypes, revealing mutually exclusive immune-active versus immune-suppressive niches, and ecotype-remodeling trajectories are traced, enhancing survival and therapy response prediction beyond standard clinical markers. SCOPE provides a more comprehensive view of tissue microenvironments and shows great potential in improving patient survival prediction and therapy response assessment.

C1011: Bayesian multi-sample and multi-scale clustering with feature selection for spatial transcriptomics data

Presenter: Alvin Sheng, University of Minnesota, Twin Cities, United States

Co-authors: Thierry Chekouo, Sandra Safo

Recent advances in spatial transcriptomics have enabled researchers to measure gene expression at the single-cell spatial resolution, often for multiple tissue sections in a single study. The aim is to present a Bayesian method that simultaneously performs factor analysis and spatial clustering on multiple samples, where the clustering is done at both the single-cell and tissue regional scale. Therefore, the method simultaneously sorts the cells into cell types and partitions the cells into spatial domains of tissue. To increase interpretability, a feature selection mechanism is employed within the estimation of the sparse factor loadings matrix, which detects genes that discriminate between cell-type clusters. The advantages of the method are illustrated over alternative state-of-the-art approaches through simulation studies and three real data applications.

C1085: Spatial Poisson-lognormal pathway model for detecting spatially expressed (SE) genes in spatial transcriptomics

Presenter: Emmanuel Sarfo Fosu, Baylor University, United States

Co-authors: Joon Jin Song, Thierry Chekouo

Spatial transcriptomics (ST) enables high-resolution mapping of gene expression across tissues, offering spatial insights into cellular organization, tissue development, disease progression, and treatment response. One key objective in ST analysis is to identify spatially expressed (SE) genes. Most existing methods, however, ignore gene-gene dependencies. The aim is to propose a spatial Poisson lognormal model that uses biological pathways to jointly capture both spatial and gene-level dependencies. Given the high dimensionality of ST data, a non-spatial conditional autoregressive (CAR) prior is adopted that models gene dependencies by borrowing external biological knowledge. The Bayesian model simultaneously detects gene clusters of non-SE and SE. By integrating these localized gene dependencies into a hierarchical spatial framework, the model improves sensitivity and interpretability in detecting SE genes. Simulation studies and applications to real ST datasets demonstrate enhanced power and accuracy compared to existing univariate methods, while leveraging biologically meaningful gene-gene relationships.

CO205 Room BCB 405 BAYESIAN STATISTICS IN BIOLOGICAL AND ENVIRONMENTAL SCIENCES

Chair: Xueying Tang

C0764: Efficient Bayesian semiparametric modeling and variable selection for spatiotemporal transmission of multiple pathogens

Presenter: Nikolay Bliznyuk, University of Florida, United States

Co-authors: Xueying Tang

Mathematical modeling of infectious diseases plays an important role in the development and evaluation of intervention plans. These plans, such as the development of vaccines, are usually pathogen-specific, but laboratory confirmation of all pathogen-specific infections is rarely available. If an epidemic is a consequence of the co-circulation of several pathogens, it is desirable to jointly model these pathogens in order to study the transmissibility of the disease to help inform public health policy. A major challenge in utilizing laboratory test data is that it is not available for every infected person. Appropriate imputation of the missing pathogen information often requires a prohibitive amount of computation. To address it, the earlier hierarchical Bayesian multi-pathogen framework is extended, which uses a latent process to link the disease counts and the lab test data. Under the proposed model, imputation of the unknown pathogen-specific cases can be effectively avoided by exploiting the relationship between multinomial and Poisson distributions. A variable selection prior is used to identify the risk factors and their proper functional form respecting the linear-nonlinear hierarchy. The efficiency gains of the proposed model and the performance of the selection priors are examined through simulation studies and on a real data case study from hand, foot, and mouth disease (HFMD) in China.

C0598: Suppressing odds ratio inflation: Detection and correction of perfect separation in logistic regression

Presenter: Liangliang Zhang, Case Western Reserve University, United States

Logistic regression is a core method for modeling binary outcomes, but it struggles when perfect separation occurs, leading to inflated odds ratios and unstable predictions. Most existing methods address this issue post hoc. In contrast, a pre-hoc linear programming approach is proposed that assesses the extent to which predictor combinations are separated by the binary outcome. When separation arises from a linear combination of multiple variables, this is defined as latent perfect separation. The method detects both direct and latent separation. It is shown that although latent separation results in infinite estimates for individual coefficients, the ratio of coefficients converges to a fixed constant that reflects the true underlying relationship. This finding is incorporated into a Bayesian power prior framework to correct inflated estimates and guide them toward realistic values. The Bayesian method also integrates additional statistical criteria to ensure convergence. Simulations show that our approach significantly outperforms the widely used Firth correction, and real-world applications confirm that our adjusted coefficients better reflect true effects. By detecting and addressing separation before model fitting, the method improves both interpretability and accuracy in logistic regression, especially in complex settings like disease association studies.

C0958: A Bayesian approach to produce subnational population estimates using a population base statistical register

Presenter: Jairo Fuquene, UC Davis, United States

Subnational population estimates (SPE) in Latin America are useful to implement new public policies in subnational areas with internal armed conflicts or difficult to access. The aim is to propose to combine a Population Base Statistical Register (PBSR) and the Official Population Projections (OPP) using a Bayesian approach to produce SPE. The proposed procedures are useful for computing the SPE of the population size or the SPE of the population size in percentage SPE (%). However, the focus is on SPE (%) due to some data restrictions and to ensure data confidentiality. The PBSR is constructed using multiple administrative sources with registers from the health, education, vital statistics systems, tax registration, and, more importantly, the registers of the victims of the current internal armed conflict in Colombia. New fast Markov chain Monte Carlo algorithms are also proposed to produce SPE (%) using data augmentation procedures to address the complications caused by the resulting joint posterior containing gamma functions. The proposal is implemented to compute SPE (%) by age and sex groups in the municipality

of Jamundi in Colombia, which is currently affected by poverty, forced displacement, and the internal armed conflict, and the accuracy is evaluated with a Population Census.

C0205: Bayesian region selection and prediction in Poisson regression with spatially dependent global-local shrinkage prior

Presenter: **Zihan Zhu**, Case Western Reserve University, United States

Co-authors: Xueying Tang, Shuang Zhou

A spatially dependent global-local shrinkage prior is proposed for Poisson regression, specifically aimed at prediction and region selection with spatially dependent covariates. This approach is inspired by the challenge of predicting the number of hurricanes and identifying regions with significant contributions based on spatially dependent data. The proposed prior combines the conditional autoregressive (CAR) prior, which introduces spatial dependence in the coefficients of spatially dependent covariates, with the super heavy-tailed (SH) prior, which ensures appropriate global-local shrinkage effects for selection. Metropolis-within-Gibbs sampler is developed for computation. Extensive simulation studies demonstrate that the method excels when signals are weak and adjacent, and the spatial dependence in covariates is strong. Applied to North Atlantic hurricane prediction, the method outperforms traditional regression-based approaches and rivals the benchmark "oracle" model.

CO026 Room BCB 406 ADVANCES IN NETWORK ANALYSIS AND NONPARAMETRIC STATISTICS

Chair: Joshua Cape

C0980: Faithful group Shapley value

Presenter: **Weijing Tang**, Carnegie Mellon University, United States

Data Shapley is an important tool for data valuation, which quantifies the contribution of individual data points to machine learning models. In practice, group-level data valuation is desirable when data providers contribute data in batches. However, it is identified that existing group-level extensions of data Shapley are vulnerable to shell company attacks, where strategic group splitting can unfairly inflate valuations. Faithful group Shapley value (FGSV) is proposed, which uniquely defends against such attacks. Building on original mathematical insights, a provably fast and accurate approximation algorithm is developed for computing FGSV. Empirical experiments demonstrate that the algorithm significantly outperforms state-of-the-art methods in computational efficiency and approximation accuracy, while ensuring faithful group-level valuation.

C1075: Fast convergence of a federated expectation-maximization algorithm

Presenter: **Rajita Chandak**, EPFL, Switzerland

Data heterogeneity has been a long-standing bottleneck in studying the convergence rates of federated learning algorithms. The benefits of data-heterogeneity are illustrated through establishing convergence rates of the expectation-maximization (EM) algorithm for the federated mixture of k -linear regressions model (FMLR). The convergence rate of the EM algorithm is completely characterized under all regimes of m/n , where m is the number of nodes and n is the number of data points per node. Furthermore, theoretical and empirical implications of various standard assumptions are discussed in the literature, showcasing the need for careful consideration of statistical frameworks for federated algorithms.

C1365: Covariate assisted graph matching

Presenter: **Jesus Arroyo**, Texas A&M University, United States

Co-authors: Trisha Dawn

Data integration is essential across diverse domains. A crucial initial step in this process involves merging multiple data sources based on matching individual records. When the datasets are network data, this problem is typically addressed through graph matching methodologies. For such cases, auxiliary features or covariates associated with nodes or edges can be instrumental in achieving improved accuracy. However, most existing graph matching techniques do not incorporate this information. To overcome these limitations, the aim is to propose two novel covariate-assisted seeded graph matching methods. The first one utilizes the quadratic assignment problem (QAP), while the second one leverages the local neighborhood structure of non-seed nodes to guide the matching process. Theoretical guarantees are established for model estimation error and exact recovery of the solution of the QAP, demonstrating perfect alignment accuracy with high probability under sufficient signal strength. The effectiveness of the methods is demonstrated through numerical experiments. Finally, the proposed approach is applied to match two real-world networks. By leveraging additional covariates such as institution, country, and graduation year, improved alignment accuracy is achieved. The power of integrating covariate information is highlighted in the classical graph matching setup, offering a practical and improved framework for combining network data with wide-ranging applications.

C1412: Perturbation-robust predictive modeling of social effects by network subspace generalized linear models

Presenter: **Can Minh Le**, University of California, Davis, United States

Co-authors: Tianxi Li, Jianxiang Wang

Network-linked data, where multivariate observations are interconnected by a network, are becoming increasingly prevalent in fields such as sociology and biology. These data often exhibit inherent noise and complex relational structures, complicating conventional modeling and statistical inference. Motivated by empirical challenges in analyzing such data sets, the purpose is to introduce a family of network subspace generalized linear models designed for analyzing noisy, network-linked data. A model inference method is proposed based on subspace-constrained maximum likelihood, which emphasizes flexibility in capturing network effects and provides a robust inference framework against network perturbations. The asymptotic distributions of the estimators are established under network perturbations, demonstrating the method's accuracy through extensive simulations involving random network models and deep-learning-based embedding algorithms. The proposed methodology is applied to a comprehensive analysis of a large-scale study on school conflicts, where it identifies significant social effects, offering meaningful and interpretable insights into student behaviors.

CO318 Room BCB 408 THEORY AND APPLICATION OF SAMPLING ALGORITHMS

Chair: Qian Qin

C0220: Metropolis-adjusted subdifferential Langevin algorithm

Presenter: **Ning Ning**, Texas A&M University, United States

The Metropolis-adjusted Langevin algorithm (MALA) is a widely used Markov chain Monte Carlo (MCMC) method for sampling from high-dimensional distributions. However, MALA relies on differentiability assumptions that restrict its applicability. The Metropolis-adjusted subdifferential Langevin algorithm (MASLA) is introduced, a generalization of MALA that extends its applicability to distributions whose log-densities are locally Lipschitz, generally non-differentiable, and non-convex. The theoretical foundation of MASLA is established by proving its convergence to a set-valued differential inclusion equation, ensuring well-defined long-run behavior. Furthermore, the performance of MASLA is evaluated by comparing it with other sampling algorithms in settings where they are applicable. Results demonstrate the effectiveness of MASLA in handling a broader class of distributions while maintaining computational efficiency.

C1242: Proximal Hamiltonian Monte Carlo

Presenter: **Eric Chi**, University of Minnesota, United States

Co-authors: Dootika Vats, Apratim Shukla

Modern statistical learning problems often face challenges in efficient sampling. This is due to a dearth of effective sampling strategies, particularly for high dimensions. In addition, this is even more pronounced for problems of image denoising and sparse signal recovery, etc. A major issue is the non-differentiability of the underlying Bayesian posterior density. This is a direct consequence of employing a non-differentiable prior, which is popular to induce sparsity in the model. As a result, sampling even from efficient gradient-based Markov chain Monte Carlo (MCMC) methods

becomes difficult. This problem is circumvented by proposing a proximal Hamiltonian Monte Carlo (p-HMC) algorithm, which uses tools like proximal mappings and Moreau-Yosida (MY) envelopes within Hamiltonian dynamics. The contribution is that conditions for geometric ergodicity of the underlying HMC chain and a methodology to obtain a suitable choice for the regularization parameter in the MY envelope are also provided. The method has been implemented for a sparse logistic regression and a low-rank matrix estimation problem, which demonstrates its efficiency over the current state of the art.

C1427: Antithetic noise in diffusion models

Presenter: Guanyang Wang, Rutgers University, United States

The purpose is to initiate a systematic study of antithetic initial noise in diffusion models. Across unconditional models trained on diverse datasets, text-conditioned latent-diffusion models, and diffusion-posterior samplers, it is found that pairing each initial noise with its negation consistently yields strongly negatively correlated samples. To explain this phenomenon, experiments and theoretical analysis are combined, leading to a symmetry conjecture that the learned score function is approximately affine antisymmetric (odd symmetry up to a constant shift), and evidence is provided supporting it. Leveraging this negative correlation, two applications are enabled: (1) enhancing image diversity in models like stable diffusion without quality loss, and (2) sharpening uncertainty quantification (e.g., up to 90% narrower confidence intervals) when estimating downstream statistics. Building on these gains, the two-point pairing is extended to a randomized quasi-Monte Carlo estimator, which further improves estimation accuracy. The framework is training-free, model-agnostic, and adds no runtime overhead.

C1483: Recurrence and transience of Markov chains and evaluation of improper priors

Presenter: Kshitij Khare, University of Florida, United States

Sufficient conditions are developed for recurrence and transience of irreducible Markov chains on the real line. The results are developed using a combination of drift (Lyapunov) conditions and increment analysis, which is based on the moments of the (one-step) jumps of the chain. The new results are applied to a simple, but fundamental problem in statistical decision theory. Specifically, suppose $U \sim N(t, 1)$ and let $a(t)$ be an improper prior density that yields a proper posterior density, $p(t|U = u)$. Consider the Markov chain on the real line whose one-step transition from x is obtained by first drawing $t \sim p(\cdot|x)$, and then adding a standard normal noise to t . Previous results in the literature imply that, if this Markov chain is (null) recurrent, then the prior $a(t)$ is strongly admissible, which basically means that the Bayes estimators generated by $p(t|u)$ are admissible. The new sufficient conditions for recurrence are used to show that various improper priors for t , including improper versions of popular shrinkage priors, are strongly admissible.

CO142 Room BCB 409 STATISTICS IN NEUROSCIENCE I

Chair: Kristin Linn

C0561: Covariate-assisted community detection for functional brain networks: A variational approach

Presenter: Panpan Zhang, Vanderbilt University Medical Center, United States

Understanding the human brain requires modeling its complex and high-dimensional architecture. Graph-based methods for brain networks, particularly community detection techniques, offer critical insights into the modular organization of brain function across different disease stages. A novel covariate-assisted community detection method is proposed that jointly leverages network topology and informative node-level attributes. The method is formulated within a variational estimation framework, enabling efficient estimation in large-scale networks. Through extensive simulation studies, the advantages of the approach are demonstrated over existing alternatives in terms of accuracy and robustness. The proposed method is further applied to fMRI data from individuals with Alzheimer's disease, uncovering meaningful community structures that reflect disease-related alterations in brain connectivity.

C0705: Modeling the neural response to head movement

Presenter: Martin Lindquist, Johns Hopkins University, United States

Head motion during functional magnetic resonance imaging (fMRI) studies can introduce both artifacts and motion-evoked neural responses, potentially biasing inferences about brain activity. Most correction methods treat motion purely as noise, overlooking its possible neural effects. Using large-scale meta-analytic brain maps, the spatial and temporal features of head-motion-evoked BOLD responses are characterized during resting-state scans in healthy young adults. A systematic BOLD increase of 3.6 seconds is observed after motion, particularly in motor-related regions of the brain. Based on this finding, FD-regression is proposed, a method that convolves thresholded framewise displacement (FD) measures with a modeled motion-evoked response to create a nuisance regressor for the GLM. FD-regression more effectively reduced motion-related signal in resting-state data compared to standard motion correction techniques. These findings suggest FD-regression is a promising approach for improving fMRI analyses.

C0759: New developments in implementing network science in brain behavior linking

Presenter: Selena Wang, Indiana University School of Medicine, United States

The aim is to propose a latent space statistical network analysis (LatentSNA) that implements network science in a generative Bayesian framework, preserves the neurologically meaningful brain topology, and improves the statistical power for imaging biomarker detection. LatentSNA (1) addresses the lack of power and inflated Type II errors in current analytic approaches when detecting imaging biomarkers, (2) allows unbiased estimation of biomarkers' influence on behavior variants, (3) quantifies the uncertainty and evaluates the likelihood of the estimated biomarker effects against chance and (4) ultimately improves brain-behavior prediction in novel samples and the clinical utilities of neuroimaging findings. LatentSNA is broadly applicable to multiple neuroimaging landmark studies, imaging modalities, and outcome measures with developing, aging, and transdiagnostic populations, totaling 8,003 to 11,861 participants. In these applications, LatentSNA achieves substantial accuracy gains (averaging 110% - 150%) and replicability improvements (averaging 153% in moderate-to-large datasets). As a result, LatentSNA provides an unprecedented view of how network topology is implicated in brain-behavior relationships.

C0943: Learning directed brain connectomes for neurodevelopment

Presenter: Ying Guo, Emory University, United States

In recent years, connectome-based research has emerged as a central focus in neuroscience, offering critical insights into brain organization and supporting predictive modeling of cognitive, behavioral, and clinical outcomes. While most existing approaches analyze undirected brain connectomes, they overlook the directionality and causal influence between brain regions. To address this limitation, the aim is to propose a novel, task-aware framework for learning directed brain connectomes from fMRI data. The method leverages advanced statistical modeling and machine learning to perform regularized inference, yielding sparse, directed connectivity graphs that capture causal interactions across the brain. Simultaneously, the method learns low-dimensional graph embeddings that are optimized to predict demographic, behavioral, and clinical outcomes. Applied to a large-scale neurodevelopmental study, the approach uncovers directed whole-brain connectivity patterns among children and adolescents and reveals new insights into subpopulation differences in the directed connectome, highlighting its potential to advance both mechanistic understanding and predictive modeling in neuroscience.

CC395 Room BCB 210 FORECASTING

Chair: Jorge Gonzalez Cazares

C1171: Reducing food waste in grocery retailing: Improved demand forecasting for specific challenges

Presenter: Maya Natascha Vienken, Bielefeld University, Germany

Co-authors: David Winkelmann, Roland Langrock, Theresa Elbracht

A key challenge for grocery retailers is to effectively manage inventories for perishables. This requires precise forecasting of customer demand to ensure customer satisfaction through a reliable service level while minimising spoilage with its financial and environmental consequences. Despite the availability of various forecasting approaches in research and practice, the dynamic nature of grocery store operations complicates accurate predictions. For instance, retailers frequently apply price reductions to units nearing expiration. Fully including such demand at a reduced price into the forecast can lead to overestimating future demand due to an uplift effect. Conversely, stockouts result in censored demand data, further complicating accurate demand forecasting. These complexities are examined using data from a major European grocery retailer, aiming to enhance demand forecasting accuracy. The ultimate objective is to contribute to the reduction of food waste in the retail sector, a commitment that major retail companies have made in partnership with the German Federal Ministry.

C1288: Short-term probabilistic energy load forecasting with GAMLSS framework

Presenter: **Katarzyna Chec**, Wrocław University of Science and Technology, Poland

Co-authors: Bartosz Uniejewski, Rafal Weron

Electricity demand forecasts are among the key determinants of electricity prices and are therefore widely used in energy price forecasting models. In practice, decision makers often rely on forecasts provided by transmission system operators (TSOs), which tend to exhibit limited accuracy and systematic bias. Recent research has shown that these forecasts can be significantly improved using simple autoregressive models with day-ahead information, yielding notable reductions in forecast errors. However, the focus is solely on point forecasts, without addressing probabilistic approaches. The aim is to fill this gap by investigating the applications of probabilistic forecasting methods for improving TSO electricity demand predictions. Both post-processing of point forecasts (i.e., building predictive distributions based on point predictions) and direct distributional forecasting are considered, with a focus on generalized additive models for location, scale, and shape (GAMLSS). The analysis covers nine years of data (2016-2024) from three major European electricity markets. Forecast accuracy is evaluated using the pinball score and statistical significance is verified with the Diebold-Mariano test. It is found that GAMLSS outperforms post-processing approaches as well as quantile regression, and the use of corrected load forecasts yields superior performance over raw TSO predictions when employed as inputs to probabilistic models.

C1345: Prediction limits based on weighted model combinations and CRPS scoring rules

Presenter: **Paolo Vidoni**, University of Udine, Italy

Co-authors: Valentina Mameli

The purpose is to define prediction limits based on predictive distributions obtained as weighted combinations of elementary distribution functions. In particular, linear and multiplicative combinations of density functions, as well as linear combinations of quantile functions, are considered. A combined model can serve as a useful surrogate for the true, unknown predictive model of the random phenomenon of interest. The weights associated with the individual components are determined by considering the continuous ranked probability score (CRPS) and its weighted extensions. Simulation studies show that, by using an appropriately weighted version of the CRPS that focuses on the tails of the distribution, the estimated combined model provides prediction limits with coverage probabilities close to the target nominal value. The good performance of the approach is further illustrated using real data on athletics' records.

C1420: A comparison of sequential and integrated estimation of realized volatility models used in portfolio optimization

Presenter: **Erik-Jan Senn**, University of St. Gallen, Switzerland

Portfolio optimization relies on precise volatility forecasts. However, in many cases, volatility models are estimated independently of the task-specific problem, leading to forecasts that may not be optimal for the intended application. The purpose is to discuss the simple example of a mean-variance utility investor seeking to maximize expected utility. Two approaches to estimating HAR-type realized volatility models are compared: i) the standard sequential procedure, which first estimates the realized volatility model based on a statistical loss function such as mean squared error, and then applies the plug-in principle to compute optimal portfolio weights and realized utilities; ii) the integrated procedure, which uses the same HAR-type model but estimates the parameters using the realized utility of the mean-variance investor as the loss function. The simulation-based and empirical comparison can guide the choice of loss functions and evaluation metrics for forecast evaluation in this setting.

CC403 Room BCB 211 APPLIED ECONOMETRICS

Chair: Juan Manuel Rodriguez-Poo

C0599: The drivers of emission reductions in the European carbon market

Presenter: **Jamie Cross**, Melbourne Business School, Australia

Co-authors: Hilde Bjorland, Felix Kapfhammer

The purpose is to study the drivers of emission reductions in the carbon market of the European Union Emission Trading System (EU ETS) since its inception in 2005. A novel empirical framework is introduced that facilitates the joint identification of simultaneous demand and supply shocks underlying the European carbon market. It is found that emission supply restrictions of the EU ETS were the dominant driver of emissions reductions, reducing emissions by 46%. However, it is also found that two opposing emission demand factors also played an important role. Demand from industrial economic activity increased emissions by 15%, while other demand-side factors, primarily reflecting the transition to low-carbon economies, reduced emissions by 21%.

C1310: Indo-Pakistani arms race dynamics: A structural VAR perspective

Presenter: **Hui Ma**, Georg-August-University Goettingen, Germany

Co-authors: Lennart Empting, Helmut Herwartz, Simon Albert Kuehn

The longstanding rivalry between India and Pakistan is regarded as a classic case of arms race dynamics, commonly analyzed with vector autoregressive (VAR) models. Previous studies often rely on Granger causality tests from reduced-form estimates to infer the direction of military buildup. Structural vector autoregressive (SVAR) models are applied to examine the effects of independent structural shocks on the military burdens of India and Pakistan from 1956 to 2022. Employing independent component analysis for data-driven identification, country-specific shocks are isolated, and their respective impacts are assessed both domestically and across borders. Findings reveal a marked asymmetry. India's military burden responds almost exclusively to its own domestic shocks, while Pakistan's military expenditure is shaped by both its own shocks and lagged responses to Indian shocks, indicating a one-sided arms race dynamic. Notably, although the immediate impact of Indian shocks on Pakistan is not statistically significant, substantial reactive effects, consistent with the Richardson models action-reaction mechanism, develop over time.

C0202: Measuring labor market precarity: A multidimensional econometric approach

Presenter: **Petros Kosmas**, Cyprus University of Technology, Cyprus

A multidimensional framework is proposed to identify and measure labor market precarity by integrating three indicators: The precariousness indicator (capturing employment instability), the vulnerability indicator (assessing economic insecurity), and the precarity indicator, which reflects their intersection. Using microdata from labor market surveys, the framework applies a structured econometric approach to (a) examine the determinants of precarious employment and economic vulnerability and (b) estimate the share of employees simultaneously affected by both conditions, thereby entering a state of precarity. This methodological design offers a systematic means of analyzing labor market fragility and its broader implications for workforce stability and economic resilience.

C1338: Assessing political, economic, and environmental drivers of fire activity in the Brazilian Amazon

Presenter: **Ilka van de Werve**, Vrije Universiteit Amsterdam, Netherlands

Co-authors: Thales Pupo West, Niek Scherpenhuijzen

The purpose is to examine the impact of political conditions on fire activity in the Brazilian Amazon, while controlling for economic and environmental circumstances. A state space model is specified with stochastic level and seasonal components using monthly data. Explanatory variables include an intervention variable for right-wing presidencies (allowing for anticipation effects), total precipitation, and growth in international commodity prices, alongside the exchange rate and inflation. Economic variables are smoothed by averaging price growth over short, medium, and long lag windows. It is found that periods of right-wing presidencies are associated with more forest fires in Brazil, a pattern not observed in control countries Peru and Colombia. Precipitation is negatively associated with fire activity. Moreover, results suggest that farmers do not respond directly to price changes but do so in the longer run, in a positive direction. Outliers and structural breaks can be explained by weather conditions and the early months of Bolsonaro's presidency.

CC441 Room BCB 407 STATISTICAL INFERENCE AND COMPUTATION
Chair: Nilanjan Chakraborty
C0376: Belief distortion under stress using trust decay and ℓ_1 -regularized inference
Presenter: **Gabriel Nixon**, New York University, United States

When agents operate under stress, they tend to overreact, ignore prior beliefs, or misallocate trust. A general framework is proposed to model this using ℓ_1 -regularized updates with entropy limits and decaying trust. The idea is to simulate how beliefs shift when information is incomplete, reliability degrades, or environments change. The method introduces a budgeted form of belief distortion where each update trades off staying close to prior assumptions with adapting to new shocks. Trust decay is modeled as a weighted penalty that gradually reduces the influence of old beliefs. Diagnostics are included to track regret, drawdowns, and entropy flow over time. Applications span any sequential decision setup under uncertainty, including finance, policy simulations, or dynamic systems. Several open problems are outlined, such as whether belief curvature can flag regime shifts or how repeated interventions distort entropy. The setup supports causal belief adjustments, mirror-descent-style learning, and stress-aware regularization that holds up even under regime breaks. It is designed to be modular, interpretable, and fast enough to use in real-time. The focus is not on perfect predictions but on understanding how agents should bend their beliefs without breaking them when operating under pressure.

C0256: Anytime valid and asymptotically efficient inference driven by predictive recursion
Presenter: **Vaidehi Dixit**, University of Missouri, United States

Co-authors: Ryan Martin

Data peeking is a common problem, and sampling till significance is obtained leads to erroneous conclusions and irreproducibility. Given this criticism, new approaches are desired. E-processes are a relatively new tool that quantify evidence against a null hypothesis, such that optional stopping is allowed, and combining evidence across multiple studies is permissible. Generally, consider comparing two classes of candidate models, where error rate control is desired at any stopping time. The focus is on the novel e-process construction that leverages the so-called predictive recursion (PR) algorithm. The resulting PRe-process gives anytime valid inference uniformly over stopping rules and is shown to be efficient in the sense that it achieves the maximal growth rate asymptotically, under the alternative relative to the mixture model being fit by PR. Specific applications of this methodology are presented, namely testing for monotonicity and log-concavity.

C1062: A penalized least squares approach to adaptive ridge regression
Presenter: **Keith Knight**, University of Toronto, Canada

A prior study showed that ridge regression estimates can be viewed as a weighted average of subset least squares estimates, with weights depending only on the design. Using this representation, the total weight that each predictor contributes to a given ridge estimate is computed, thereby defining the notion of partial degrees of freedom for each predictor; the equivalent degrees of freedom are then equal to the sum of the partial degrees of freedom. Thus, by varying the ridge parameters, can define "fractional" model selection where the partial degrees of freedom for each predictor vary between 0 and 1. A penalized least squares approach is considered for estimating the partial degrees of freedom, allowing the estimated partial degrees of freedom to potentially equal 0. Some asymptotic theory of the resulting estimates is also considered.

C0948: Kolmogorov-Arnold networks for high-dimensional estimation: A method of sieves approach
Presenter: **Sami Abdurahman**, Toronto Metropolitan University, Canada

The aim is to introduce a **novel sieve extremum estimator based on Kolmogorov-Arnold networks (KANs)**, designed for nonparametric estimation in high-dimensional settings with time-series data. By integrating KANs with a method of sieves estimation approach and leveraging sparsity, the method achieves asymptotic convergence rates in the $L_2(\mu)$ norm that are **independent of the covariate dimension**, thereby directly circumventing the curse of dimensionality. Specifically, the approach yields an explicit convergence rate of $op(n^{-1/4})$. This framework accommodates diverse applications such as high-dimensional conditional density estimation and nuisance function estimation within double/debiased machine learning, areas where traditional deep neural networks often struggle due to their "black box" nature, reliance on i.i.d. assumptions, and slower convergence rates. The proposed KAN sieve estimator overcomes these limitations by learning activation functions using B-splines, and its theoretical framework rigorously permits stationary mixing processes. The efficacy of this framework is demonstrated in an empirical application to high-dimensional financial time series, including predicting asset returns and volatility based on a large set of economic and market indicators.

CV413 Room BCB 207 FINANCIAL ECONOMETRICS (VIRTUAL)
Chair: Luigi Grossi
C0696: Selection confidence sets for equally-weighted portfolios
Presenter: **Sandra Paterlini**, University of Trento, Italy

Co-authors: Alessandro Fulci, Davide Ferrari

Equally-weighted portfolios offer a straightforward investment rule that can deliver strong out-of-sample performance compared with optimized portfolios. Nevertheless, the uncertainty over which equally-weighted portfolio is truly optimal among all possible combinations is typically overlooked. The selection confidence set is introduced: The collection of equally-weighted portfolios that, at a chosen confidence level and under a specified loss function, are statistically indistinguishable from the unknown optimal portfolio. The selection confidence set challenges the conventional notion of a single best portfolio. Analogously to a traditional confidence set in parameter estimation, its cardinality reflects the degree of portfolio selection uncertainty: In the presence of noisy data, distinguishing among competing portfolios becomes difficult, resulting in a larger set. Using Monte Carlo experiments with different initial assumptions and empirical data sets, it is shown that the selection confidence set is a practical tool for investors who wish to assess the robustness of their allocation choices.

C1167: ResNets and random matrix denoising on complex covariance structures: A cryptocurrency portfolio application
Presenter: **Andres Garcia**, Autonomous University of Baja California, Mexico

Covariance matrices estimated from short, noisy, non-Gaussian financial time series (particularly cryptocurrencies) are notoriously unstable. Sampling noise, heavy tails, asynchronous trading, and regime shifts inflate eigenvalues, distort principal directions, and degrade downstream portfolio optimization (e.g., Markowitz, risk parity, min-variance). Random matrix theory (RMT) offers theoretical eigenvalue estimation on high-dimensional settings, while residual neural networks (ResNets) can learn flexible, data-driven shrinkage and structure. The aim is to develop a hybrid RMT-informed ResNet denoiser for complex covariance structures and evaluate it on live-like crypto portfolio tasks.

C1255: Unstable gains: Evaluating the reproducibility of deep reinforcement learning in trading and portfolio management
Presenter: **Przemyslaw Gradzki**, University of Warsaw, Poland

The aim is to investigate the reproducibility and robustness of deep reinforcement learning (DRL) in financial applications, focusing on algorithmic trading and portfolio management across two asset classes: Stocks and cryptocurrencies. While DRL has gained popularity in these domains, most studies rely on single-run evaluations and overlook the high variance inherent to these methods. Influential DRL-based strategies are reproduced under identical hyperparameters but across multiple independent random seeds, and it is shown that both performance and learned policies vary widely under fixed configurations. These experiments highlight the fragility of commonly reported results. Even the best-performing algorithms display substantial variability across runs. To improve reliability, a checkpointing strategy is introduced, and uncertainty is quantified using bootstrapping and permutation tests. Findings reveal that prevailing evaluation practices risk misleading conclusions about strategy efficacy and also conceal the true risk profile of DRL-based financial models. This underscores the need for more rigorous and reproducible protocols to ensure dependable advancements and foster genuine risk assessment in financial DRL research

C1185: High-frequency option predictability

Presenter: **Sebastian Egebjerg**, Aarhus University, Denmark

Co-authors: Christine Bangsgaard

The purpose is to study the high-frequency predictability of S&P 500 index option returns using trade and quote data combined with machine learning methods. The models achieve substantial out-of-sample predictive power, with median R^2 around 20% and directional accuracy above 53% across the sample period. Predictability is concentrated in out-of-the-money and less liquid options, is stronger toward the end of the trading day, and quickly disappears beyond a few minutes. The most important predictors include short-term past returns and implied volatility, alongside microstructure variables such as bid-ask spreads and order-flow imbalances. Incorporating information from other parts of the option surface does not meaningfully improve forecasts, and retail trading activity is associated with reduced predictability. Overall, the results document consistent short-term predictability in SPX options, driven by market microstructure rather than systematic mispricing, and unlikely to yield exploitable trading strategies.

Authors Index

- Abba, M., 176
 Abbas, Y., 10
 Abdurahman, S., 195
 Abuqrais, M., 143
 Acharyya, A., 58
 Acharyya, S., 190
 Acquah, T., 134
 Adam, T., 136
 Adekpedjou, A., 39
 Adhikari, S., 56
 Adland, R., 180
 Adrian, T., 3
 Agarwal, A., 166
 Agiropoulos, C., 92, 180
 Agostinelli, C., 15
 Agosto, A., 88
 Agrawal, S., 150, 189
 Agterberg, J., 106
 Aguilera-Morillo, M., 79
 Ahelegbey, D., 68
 Ahipasaoglu, S., 33
 Ahmed, E., 93
 Aikman, D., 3
 Akama, Y., 10
 Akyildiz, D., 175
 Al Sadoon, M., 26
 Alaimo Di Loro, P., 124
 Alam, E., 164
 Alamichel, L., 56
 Aldieri, L., 62
 Aletti, G., 87, 88, 123
 Alexander, B., 184
 Alexander-Bloch, A., 142
 Alfiero, M., 95
 Alfo, M., 146
 Alfonzetti, G., 117
 Algeri, S., 165
 Ali, A., 61
 Alkhoury, S., 113
 Allen, S., 168
 Allison, J., 145
 Alonso, I., 85
 Alqarni, A., 34
 Alshahrani, N., 80
 Alshamari, H., 92
 Altieri, L., 52
 Altmeyer, R., 8
 Amama Ben Hassun, N., 182
 Amburgey, A., 163
 Amendola, A., 62, 108
 Amengual, D., 115
 Amezouwu, K., 107
 Amjadian, H., 22
 Ammann, K., 136
 Amongero, M., 154
 Anceschi, N., 155
 Anderson, B., 179
 Andersson, J., 180
 Andersson, V., 141
 Ando, T., 64
 Andre, J., 27, 45
 Andrei, A., 128
 Androulakis, E., 71
 Antonelli, J., 172
 Antonicelli, M., 53
 Anyaso-Samuel, S., 139
 Anzai, T., 20, 78
 Aoshima, M., 7
 Apicella, G., 99
 Aranda, S., 78
 Arbour, D., 190
 Arcagni, A., 66
 Archakov, I., 29
 Arduini, T., 183
 Arena, G., 148
 Arenas, C., 78
 Arendarczyk, M., 6
 Arezzo, M., 66
 Argiento, R., 107, 127
 Arias-Castro, E., 105
 Armann, I., 182
 Arnone, E., 69
 Arroyo, J., 192
 Arteaga Molina, L., 2
 Artemiou, A., 34
 Arya, S., 189
 Ascari, G., 99
 Ascari, R., 16, 153
 Ascorbebeitia, J., 92
 Asin, J., 72
 Aslam, N., 13
 Assenmacher, K., 157
 Assmann, L., 180
 Asti, G., 68
 Aston, J., 23
 Atak, A., 3
 Athreya, A., 106, 174
 Aubray, J., 129
 Audrino, F., 108
 Auerbach, E., 183
 Ausin, C., 72
 Avila Marquez, M., 168
 Awaya, N., 9
 Ayres, L., 33
 Azad, S., 129
 Azeraf, E., 130
 B M Souza, A., 85
 Babii, A., 45
 Bacri, T., 147
 Badics, M., 160
 Baer, B., 138
 Baghfalaki, T., 20
 Bahl, U., 183
 Bai, Y., 89
 Bailey, N., 5
 Bailey, R., 124
 Baio, G., 90
 Baione, F., 6
 Balakrishnan, S., 138
 Balasubramanian, K., 97
 Balcerek, M., 6
 Baldi Antognini, A., 87
 Ballante, E., 15
 Balocchi, C., 36
 Baltodano Lopez, O., 42
 Banbura, M., 63
 Bancroft, E., 182
 Bandeira, T., 161
 Banerjee, A., 149, 165
 Banerjee, P., 142
 Banerjee, S., 169
 Bangsgaard, C., 196
 Barbaglia, L., 45
 Barbera, J., 92
 Barbiero, A., 161
 Barigozzi, M., 116
 Barone, R., 38
 Bartlett, T., 74
 Basante, G., 115
 Basu, S., 187
 Basu, T., 184
 Bate, S., 171
 Batra, A., 123
 Battauz, M., 131
 Bauer, L., 48
 Baumeister, C., 3
 Baviera, R., 14
 Beaulieu, M., 135
 Bechlioulis, A., 180
 Beck, E., 43
 Beck, K., 180
 Beckmann, J., 85
 Bedendo, M., 11
 Bedowska-Sojka, B., 63
 Behluli, R., 43
 Behrouzi, P., 154
 Belalia, M., 112
 Bell, C., 125
 Bellio, R., 39, 117, 131
 Bellocca, G., 111
 Belton, R., 187
 Beltrami, F., 68
 Beltran-Sanchez, M., 39
 Benedetti, I., 110, 144
 Benevento, A., 134
 Benkeser, D., 17
 Bennett, J., 96
 Benoit, S., 92
 Bensoussan, A., 156
 Beranger, B., 11
 Berentsen, G., 164
 Berg, S., 188
 Berlinghieri, R., 155
 Bernardi, M., 154
 Berrettini, M., 16
 Bertagnolli, G., 117
 Bertazzi, A., 9
 Bertolacci, M., 73
 Besbeas, T., 20
 Beskos, A., 67, 105, 175
 Bessec, M., 45
 Bethlehem, R., 142
 Beutner, E., 47
 Bevington, R., 12
 Beyhum, J., 24
 Bhaduri, M., 116
 Bhaskaran, A., 47, 71
 Bhattacharjee, A., 82
 Bhattacharjee, S., 105
 Bhattacharya, A., 57
 Bhattacharya, B., 58
 Bhattacharya, I., 175
 Bhattacharya, R., 58
 Bi, J., 68
 Bia, M., 143
 Bian, J., 123
 Bianchi, D., 12, 85
 Bidder, R., 3
 Bideau, B., 74
 Bien-Barkowska, K., 21
 Bierkens, J., 150
 Bignozzi, E., 103
 Bilinska, K., 159
 Billi, R., 159
 Billio, M., 42, 120
 Birbil, I., 105, 153
 Biswas, S., 189
 Bjornland, H., 194
 Bladt, M., 10, 64, 150
 Blaise, K., 66
 Blangiardo, M., 53
 Blasques, F., 46
 Bliznyuk, N., 191
 Blotvogel, R., 49
 Bo, S., 97
 Bo, X., 126
 Bobbia, B., 60
 Bobeica, E., 63
 Bocquet-Nouaille, L., 60
 Boehning, D., 124
 Boetti, C., 15, 74
 Bonaccolto, G., 62
 Bonam, D., 99
 Bongiorno, E., 70
 Boniece, B., 133
 Bornstein, A., 127
 Borri, N., 62
 Borrotti, M., 34
 Borsch, M., 178
 Bortolato, E., 140
 Boschi, T., 52
 Botha, A., 176
 Bouchaud, J., 109
 Boudt, K., 48
 Bouveret, A., 48
 Bracher, J., 168
 Bradley, J., 175
 Braekers, R., 96
 Braun, D., 189
 Brautigam, M., 157
 Brave, S., 27
 Breitung, J., 95
 Brenner, L., 165
 Brombin, C., 14
 Brown, A., 125
 Browne, R., 20
 Brownlee, C., 85
 Bruha, J., 19
 Bruinsma, M., 8
 Brunetti, C., 185
 Brunetti, M., 62
 Bruni, B., 141
 Bruseghini, A., 68
 Bryzgalova, S., 148
 Bu, F., 190
 Buccheri, G., 84
 Bucci, A., 46
 Bucci, D., 158
 Buch, D., 36

- Buchner, H., 50
 Buckmann, M., 177
 Buizza, A., 68
 Bulla, J., 164
 Burke, K., 19, 77
 Burnello, J., 168
 Buscaglia, J., 166
 Busshoff, H., 143
 Butters, A., 27
 Bystrov, V., 12, 108

 Caballero-Aguila, R., 79
 Caeiro, F., 120
 Cai, J., 54, 167
 Cai, M., 160
 Cai, X., 58
 Cai, Z., 118, 174
 Calderazzo, S., 55
 Calissano, A., 70
 Callaway, B., 66
 Camehl, A., 121
 Campagna, A., 68
 Campos Martins, S., 95
 Candila, V., 46, 62
 Cantwell, R., 23
 Cao, J., 151
 Cape, J., 106, 141
 Capezza, C., 23
 Capitoli, G., 15, 163
 Caporin, M., 149, 150, 179
 Cappelli, C., 61
 Cappozzo, A., 14
 Caraianni, P., 150
 Carallo, G., 42
 Carbonari, L., 185
 Cardoso, D., 13
 Carfora, A., 146
 Carl, D., 11, 151
 Carmona Benitez, R., 144
 Carnero, A., 24
 Carrion-i-Silvestre, J., 149
 Casa, A., 117, 133
 Casarin, R., 42, 43, 68, 83
 Cascaldi-Garcia, D., 99
 Casoli, E., 53
 Castellanos Nueda, M., 140
 Castiglione, C., 154
 Castillo-Mateo, J., 53, 72
 Castle, J., 116
 Castro-Camilo, D., 72
 Catania, L., 135
 Caterini, G., 27
 Cattelan, M., 39
 Cavaliere, G., 29
 Cavicchioli, M., 46
 Cebrian, A., 72
 Ceccarelli, E., 53
 Cecere Palazzo, K., 46
 Cen, Z., 136
 Cepni, O., 46
 Cerchiello, P., 88
 Cerqueti, R., 65
 Chakrabarty, S., 58
 Chakraborty, N., 57, 84
 Chakraborty, S., 188
 Chambaz, A., 17
 Chan, K., 52, 70
 Chan, S., 3, 176, 177
 Chan, W., 176
 Chanatasig, E., 159
 Chance, E., 20
 Chanda, A., 165
 Chandak, R., 192
 Chandna, S., 74
 Chang, C., 126
 Chang, H., 121
 Chang, J., 41
 Chantarasap, P., 147
 Chapagain, N., 97
 Charemza, W., 85
 Chatterjee, K., 71
 Chaumeny, Y., 57
 Chavez Martinez, G., 136
 Chec, K., 194
 Chekouo, T., 32, 191
 Chen, C., 4
 Chen, D., 123
 Chen, G., 60
 Chen, H., 149
 Chen, J., 58, 180
 Chen, K., 126, 188
 Chen, L., 132
 Chen, P., 4, 98
 Chen, R., 179
 Chen, S., 128, 173
 Chen, T., 106, 174
 Chen, W., 170
 Chen, X., 97
 Chen, Y., 25, 66, 67, 117, 123, 137, 158
 Chen, Z., 123
 Cheng, T., 4
 Cheng, Y., 152
 Chenouri, S., 103
 Chevallier, J., 27
 Chi, E., 192
 Chiaromonte, F., 52
 Chib, S., 162
 Cho, Y., 152
 Choi, S., 153
 Choi, T., 153
 Chopin, N., 9
 Christev, A., 90
 Christiaens, L., 134
 Chu, A., 3
 Chu, B., 174
 Chu, J., 176, 177
 Chudik, A., 5
 Ciarreta, A., 159
 Cibersam, G., 78
 Cincinelli, P., 153
 Cinelli, C., 162
 Cipolla, S., 33
 Cipollini, A., 12, 63
 Cipriani, M., 146
 Ciuperca, G., 81
 Clark, D., 59
 Clarke, C., 127
 Clemente, G., 5
 Coffman, D., 58
 Cohen, G., 24
 Cojoianu, T., 68
 Colombo, P., 72
 Comola, M., 109
 Comunello, C., 109
 Conde-Amboage, M., 161
 Confalonieri, G., 28
 Conrad, T., 152
 Consiglio, A., 62
 Cook, A., 33
 Coolen, F., 92
 Coolen-Maturi, T., 92
 Cooley, D., 137
 Cooper, A., 61
 Coppier, R., 65
 Corberan-Vallet, A., 39
 Corenflos, A., 175
 Cormand, B., 78
 Cornea-Madeira, A., 83
 Corradin, F., 42
 Corradin, R., 155
 Correia Virtuoso Sebastiao, H., 110
 Corrigan, N., 124
 Corsaro, S., 89
 Corsi, F., 84
 Cortes Martinez, J., 182
 Cortese, F., 69
 Cory-Wright, R., 54
 cosimopetracchi, C., 185
 Cosma, G., 89
 Cossette, H., 47, 100
 Costa Andreu, E., 24
 Costa, E., 70
 Costa, M., 18
 Costantini, M., 83
 Costola, M., 120
 Craiu, R., 42
 Cremaschi, A., 17, 36, 190
 Cremers, J., 48
 Cremona, M., 52, 135
 Crespi, A., 68
 Crimaldi, I., 87, 88
 Crommen, G., 96
 Cross, J., 194
 Croux, C., 41
 Csapai, A., 43
 Cuba, D., 72
 Cubadda, G., 134
 Cuevas, F., 170
 Cui, Y., 37, 89, 170, 171
 Cuker, A., 123
 Cuparic, M., 81
 Czolkova, A., 2
 Czosnyka, M., 6

 D Ambrosio, A., 70
 D Angelo, L., 188
 da Cunha Godoy, L., 137
 Dahl, B., 181
 Dai, D., 21
 Dai, F., 51
 Dai, R., 60
 Dai, X., 103
 Dal Torrone, E., 38, 143, 183
 Dalavai, U., 74
 Dalla, V., 135
 Damico, G., 13
 DAmico, S., 68
 Dan Gabriel, A., 150
 Dancik, G., 51
 Dang, K., 73
 Dang, S., 20, 110
 Dang, U., 110
 Daniele, M., 45
 Daouia, A., 10, 151
 Darvishi, R., 32
 Das, K., 93
 Das, P., 59, 188
 Das, R., 55
 Das, S., 165
 Daskalakis, C., 37
 Dassios, A., 80
 Datta, J., 169, 181
 Dattner, I., 177
 Dauber, M., 151
 Daudel, K., 9
 Davis, E., 160
 Davison, A., 57
 Dawabsha, M., 77
 Dawn, T., 192
 Dayta, D., 116
 De Alwis, T., 116
 De Angelis, D., 15
 De Blasi, P., 154
 de Carvalho, M., 30, 90, 108
 De Diego Avila, F., 156
 De Giuli, M., 87
 de Haan, L., 102
 De Iorio, M., 17, 155
 de la Pena, V., 49
 De Luca, G., 48, 134
 de Punder, R., 168
 De Truchis, G., 46, 167
 Dean, N., 17
 Deaner, B., 119
 Deb, S., 169, 184
 Deek, R., 169
 Dees, S., 66
 Deistler, M., 179
 Dejean, S., 107
 del Barrio Castro, T., 167
 del Barrio, E., 21
 Del Tatto, V., 87
 Della Porta, M., 68
 Delleani, M., 68
 Delsol, L., 52
 Demetrescu, M., 83, 95
 DeMiguel, V., 101
 Deng, L., 47
 Denti, F., 20
 Derumigny, A., 156
 Descallar, J., 71
 Dette, H., 7
 Deutscher, C., 74
 Devijver, E., 69
 Dewaskar, M., 36
 Dey, D., 132, 140
 Dey, K., 188
 Dey, T., 188
 Di Brisco, A., 153
 Di Giorgio, G., 62
 Di Iorio, A., 89
 Di Iorio, J., 52
 Di Lorenzo, E., 99
 di Lorenzo, E., 5
 Di Sauro, V., 89

- Di Stefano, R., 66
 Diakonova, M., 85
 Dias, A., 91
 Diaz Coto, S., 93
 Diaz, I., 17
 Dillschneider, Y., 102
 Dimitriadis, T., 168
 Dissanayake, G., 63
 Ditzhaus, M., 96
 Dixit, V., 195
 Djegou, E., 39
 Dobler, D., 96
 Dobronyi, C., 114
 Doernemann, N., 118
 Dolkar, T., 170
 Dominguez Diaz, L., 168
 Dominici, F., 17
 Dong, Y., 125
 Doohan, J., 77
 Doornik, J., 116
 Dorman, K., 51
 Doz, C., 45, 95
 Dragun, K., 48
 Drovandi, C., 73
 Drumia, P., 91
 Duan, F., 44
 Duan, R., 123
 Dubiel-Teleszynski, T., 157
 Dubin, J., 172
 Ducci, N., 15
 Dufour, A., 120
 Dufour, J., 135, 179
 DUggento, A., 53
 Dumitrescu, E., 45
 Dunker, F., 21
 Dunlavy, D., 32
 Dunn-Sigouin, E., 164
 Dunson, D., 36, 155, 181
 Duran-Martin, G., 90
 Durante, D., 56, 154
 Durante, F., 134
 Durban, M., 79
 Dutta, S., 18
 Duttilo, P., 165
 Dzuverovic, E., 84
 Economidou, C., 180
 Eeles, R., 182
 Egashira, K., 7
 Egebjerg, S., 196
 Egidi, L., 75
 Eguchi, M., 159
 Elbracht, T., 193
 Empting, L., 28, 132, 194
 Engelke, S., 120
 Erlwein-Sayer, C., 113
 Ertefaie, A., 17
 Escanciano, J., 2
 Esquivel, F., 112
 Esteves dos Santos, J., 110
 Evangelaras, H., 71
 Evangelou, M., 182
 Ewald, C., 166
 Fabre, T., 101
 Fabretti, A., 62
 Faes, C., 39
 Fagerberg, G., 129
 Faias, J., 95
 Failli, D., 136, 162
 Fair, D., 142
 Fan, J., 185
 Fan, X., 111
 Fang, F., 156
 Fang, Q., 126
 Farcomeni, A., 15, 185
 Farne, M., 118
 Fassino, C., 156
 Fatouh, M., 28
 Favero, C., 28
 Fayaz, M., 144
 Febrero-Bande, M., 158
 Fei, W., 174
 Feleppa, D., 99
 Feller, A., 190
 Fellingham, G., 181
 Feng, B., 176
 Feng, H., 25
 Feng, Z., 61, 177
 Fernandez, D., 182
 Fernandez, T., 39
 Fernandez-Perez, A., 11
 Ferrara, L., 45, 95, 180
 Ferrari, D., 117, 133, 195
 Ferrari, S., 76
 Ferreira Batista Martins, I., 86
 Ferreira Martins, A., 110
 Ferreira, A., 102
 Ferreira, M., 170
 Ferrer, C., 30, 108
 Ferrigno, S., 79
 Fianu, E., 68
 Figini, S., 15
 Figueres, J., 157
 Figuerola-Ferretti, I., 166
 Finnan, D., 65
 Finocchio, G., 151
 Fiorentini, G., 115
 Fiori, A., 153
 Firpo, S., 30
 Fischer, A., 37
 Fischer, C., 68
 Fisher, A., 140
 Fitzgerald, K., 59
 Fleissner, M., 87
 Florens, J., 96
 Flournoy, N., 123
 Flowers, C., 172
 Fonseca, G., 130
 Fontana, M., 70
 Fop, M., 127
 Forastiere, L., 38, 143, 183
 Forcina Barrero, A., 68
 Forcina, D., 23
 Forni, M., 27
 Foroni, C., 157
 Fortin, I., 83
 Foss, A., 70
 Fosten, J., 158
 Fox, S., 166
 Frale, C., 27
 Francis, N., 163
 Francq, C., 76, 85, 95
 Frangiamore, F., 165
 Franus, T., 62
 Franzolini, B., 36, 190
 Freire, G., 85
 Fresoli, D., 135
 Frey, C., 86
 Friel, N., 148, 154
 Frieri, R., 87, 92
 Fries, S., 46, 167
 Fritsch, M., 132
 Fritz, C., 37
 Fu, H., 60
 Fuchs-Kreiss, A., 37
 Fuentes-Garcia, R., 64
 Fuentes-Martinez, R., 24
 Fuertes, A., 11, 12
 Fujimori, K., 7
 Fukuchi, J., 102
 Fulci, A., 195
 Fuquene, J., 191
 Furceri, D., 165
 Furno, F., 27, 121
 Furrer, R., 78
 Gadea, L., 149
 Gaffi, F., 56
 Gagneur, J., 174
 Gaigall, D., 145
 Galanos, G., 92
 Galeano, P., 24, 158
 Galimberti, G., 16
 Gallo, G., 46
 Galvao, A., 30, 96
 Gambetti, L., 27
 Ganjali, M., 20
 Gao, C., 36, 124
 Gao, J., 4
 Gao, L., 100
 Gao, P., 160
 Gao, R., 121
 Gao, S., 26
 Gao, W., 141
 Gao, Y., 89
 Gao, Z., 174
 Garcia Camacha Gutierrez, I., 70
 Garcia Sanz, A., 31
 Garcia, A., 195
 Garcia-Donato, G., 140
 Garcia-Escudero, L., 68, 76
 Garcia-Jorcano, L., 24, 31
 Garcia-Perez, A., 76
 Garcin, M., 7
 Garron Vedia, I., 111
 Gatu, C., 91
 Gaudlitz, S., 8
 Ge, S., 4
 Gelain, P., 157
 Gelein, B., 74, 107
 Gelfand, A., 72
 Gelinas, A., 47
 Gerlach, R., 4
 Gerontogianni, L., 74
 Gerstenberg, J., 145
 Ghahramanpour, S., 28
 Ghani, D., 12
 Ghashti, J., 110
 Ghassami, A., 97
 Ghebremichael, M., 54
 Gherardini, L., 154
 Ghidini, V., 127
 Ghiglietti, A., 87, 88
 Ghosal, R., 104
 Ghosh, D., 84, 164, 175
 Ghosh, S., 84
 Ghoudi, K., 86
 Ghysels, E., 24, 45
 Giammaria, A., 63
 Giampino, A., 16
 Giannerini, S., 179
 Giannone, D., 3, 27, 100, 121, 165
 Giddings, K., 110
 Gigliarano, C., 89
 Giglio, C., 157
 Gilbert, P., 38
 Gilmour, S., 71, 124, 171
 Giordano, S., 63
 Giovannelli, A., 95
 Giraitis, L., 29, 135
 Girolimetto, D., 150
 Giudici, P., 15
 Giummole, F., 130
 Glargaard, L., 150
 Glielmo, A., 87
 Gloter, A., 150
 Gnecco, N., 120
 Goegebeur, Y., 10
 Goessler, G., 21
 Goessler, W., 21
 Goia, A., 52, 70
 Gomes, I., 120
 Goncalves de Souza, E., 143
 Goncalves, S., 27, 48
 Gondzio, J., 33
 Gonzalez Cazares, J., 64
 Gonzalez-Manteiga, W., 158, 161
 Gonzalo, J., 149
 Goo, J., 144
 Goodwin, T., 47
 Goracci, G., 179
 Gorji, M., 164
 Gottard, A., 18, 117
 Goulby, Z., 45
 Gourieroux, C., 30
 Goyal, V., 148
 Gracia, Z., 72
 Gradzki, P., 195
 Graham, M., 175
 Grassi, A., 157
 Grassi, S., 95
 Greenwood-Nimmo, M., 64
 Grevén, S., 2
 Griffin, J., 1, 64, 127
 Grilli, L., 15
 Grith, M., 151
 Grivas, C., 92
 Grondelli, M., 68
 Grossi, L., 68
 Grossmann, H., 123
 Gruen, B., 12
 Grzesiek, A., 6
 Gu, J., 114

- Gu, Y., 57
Gu, Z., 26
Guagnano, G., 66
Gudziunaite, S., 53
Guerrieri, L., 100
Guevara, I., 117
Guglielmi, A., 38
Guha Niyogi, P., 59
Guha, A., 57
Guha, N., 169
Guha, R., 122
Guhathakurta, K., 168
Gui, T., 138
Guidotti, E., 101
Guillaumin, A., 81
Guillotel, A., 74
Guillou, A., 10
Guindani, M., 127, 140
Gundersen, K., 146, 164
Guo, H., 69
Guo, S., 126
Guo, Y., 193
Gupta, A., 109
Gupta, K., 184
Gupta, M., 189
Gupta, R., 46
Gutierrez, L., 117
Gyger, T., 78
Gyotoku, Y., 150
- Hacquard, O., 152
Hadj-Amar, B., 127
Haefke, C., 43
Haimerl, P., 88
Halbleib, R., 29
Hallam, M., 115
Hallin, M., 1
Hambuckers, J., 119, 134
Hamid, J., 173
Hamilton, K., 123
Han, D., 8
Han, L., 36
Han, X., 101
Han, Y., 137
Hanck, C., 48, 168
Hanka, J., 166
Hans-Peter, P., 73
Hansen, S., 140
Hao, J., 111
Happersberger, D., 86
Hara, H., 160
Hara, N., 75
Hardcastle, L., 90
Harezlak, J., 128
Harhay, M., 35
Harris, J., 185
Harshaw, C., 37
Harvey, A., 135, 145
Hasanov, F., 184
Haslbeck, J., 148
Hattori, S., 54
Haupt, H., 132
Hayashi, T., 67
Hayes, A., 185
Hazelton, M., 59
He, F., 108
He, S., 49
- He, Y., 173
Heard, N., 12
Hecq, A., 134
Hector, E., 141
Heinekamp, S., 178
Heiner, M., 181
Heinrichs, F., 19
Henderson, D., 3, 132, 142, 167
Hendler, M., 6
Hendry, D., 116
Hennig, C., 70
Henningsen, T., 150
Herath, W., 187
Hermes, S., 154
Hernandez, P., 79
Herrera, A., 27
Herrmann, K., 134
Herwartz, H., 28, 157, 194
Hettinger, G., 56
Hickman, M., 15
Higdon, D., 178
Hill, E., 19
Hiraki, D., 8
Hirsch, S., 168
Hirt, C., 156
Hirukawa, J., 7
Hlouskova, J., 83
Hlubinka, D., 23
Ho, N., 57, 73
Hoelleland, S., 147, 164
Hoermann, S., 159
Hofer, V., 21
Hoffman, K., 17
Hofmarcher, P., 12, 13
Holland, C., 131
Holling, H., 123
Honaker, M., 121
Honda, T., 76
Hong, D., 106
Hong, E., 91
Hong, S., 80
Hornung, R., 50, 51
Horvath, L., 133
Hossain, S., 33
Hounyo, U., 30
Hovhannisyan, A., 108
Hrachov, M., 72
Hristopulos, D., 81
Hron, K., 2
Hsiao, C., 158
Hsu, Y., 4
Hu, G., 151
Hu, J., 79
Hu, L., 56
Hu, Y., 97, 189
Hu, Z., 132
Hua, W., 190
Huang, C., 88
Huang, W., 141
Huang, X., 172
Huang, Y., 58
Huang, Z., 117
Hubbard, R., 123
Hubner, P., 119
Hubner, S., 119
Hunter, D., 37
- Huth, K., 148
Hyodo, M., 88
- Iacopini, M., 86
Iacoviello, M., 100
Iafrate, F., 151
Iania, L., 49, 121
Iascone, E., 68
Ibragimov, R., 50
Iguchi, Y., 67
Ikeda, K., 116
Ilangasekara, S., 181
Imaizumi, M., 105
Inacio, V., 117
Inan, E., 134
Ince, A., 22
Inose, S., 17
Irigoien, I., 78
Isenberg, D., 35
Isler, C., 130
- Jablonski, I., 6
Jackson Young, L., 115
Jaffrezic, F., 117
Jaganjac, D., 186
Jagannath, A., 103
Jahan-Parvar, M., 185
Jakovac, A., 13
Jalilian, A., 170
Jana, S., 185
Janczura, J., 7
Jang, P., 122
Jara, A., 181
Jasiak, J., 48
Jeffrey, E., 81
Jenkins, P., 9, 38
Jessup, S., 100
Jewson, J., 153
Jiang, J., 133
Jiang, P., 104, 115
Jiayu, C., 138
Jimenez, I., 24
Jimenez-Martin, J., 31
Jin, C., 25
Jin, J., 26
Jin, Y., 53
Jin, Z., 153
Johansen, A., 175
Johnson, B., 174
Johnstone, J., 33
Jona Lasinio, G., 53, 72
Jones, D., 44
Jones, M., 142
Jongbloed, G., 8
Jorda, V., 2
Ju, N., 125
Jung, R., 111
Jung, S., 117
Juodis, A., 44
- Kaino, Y., 93, 102
Kalaria, S., 183
Kalli, M., 127
Kalogeropoulos, K., 157
Kaltenbach, H., 171
Kamatani, K., 150
Kandiros, V., 37
Kanellopoulos, C., 92
- Kanfer, F., 128, 188
Kang, K., 142
Kang, S., 103, 188
Kano, Y., 160
Kao, M., 131
Kapetanios, G., 29, 45, 135
Kaphammer, F., 194
Karabatsos, G., 129
Karadimitropoulou, A., 180
Karantali, M., 109
Kareem, K., 51
Karlsson, P., 164, 181
Karouzakis, N., 157
Kasprowicz, M., 6
Kastner, G., 17
Kawamo, E., 150
Kawasaki, Y., 10
Kayal, T., 129
Ke, T., 26, 113
Keeling, M., 160
Keller, A., 142
Kelner, M., 177
Kennedy, E., 52, 138
Keresztely, T., 13
Keribin, C., 188
Kessler, D., 185
Khalaf, L., 135
Khalili, A., 32, 136
Khare, K., 193
Khorrami Chokami, A., 153
Kiermeier, M., 90
Kilian, L., 27
Kim, C., 121
Kim, D., 119
Kim, J., 50, 98, 175
Kim, K., 34, 114
Kim, M., 97
Kim, S., 148, 152
Kim, Y., 19
Kimmel, S., 123
Kindji, G., 182
Kirch, C., 62
Kiriliouk, A., 11
Kirk, P., 16, 57
Kleynhans, A., 188
Klug, S., 130
Knight, K., 195
Knight, M., 15, 74
Knipp, C., 185
Knotek, E., 142
Knueppel, M., 75
Ko, S., 60
Kobayashi, G., 8
Koh, J., 31, 48
Kohlschmidt, J., 55
Kohn, R., 129
Koike, Y., 67
Kojadinovic, I., 41
Kojima, S., 4
Kokoszka, P., 118, 158
Kolaczyk, E., 173
Kolar, M., 26
Komendantova, N., 22
Kon Kam King, G., 38
Kontoghiorghe, E., 91
Koo, B., 5
Koohi-Kamali, F., 99

- Koop, G., 26, 45
 Koopman, S., 46
 Korajczyk, R., 149
 Korner, H., 19
 Koskela, J., 9
 Koslovsky, M., 181
 Kosmas, P., 194
 Kosorok, M., 54
 Kossovsky, A., 65
 Kote-Jarai, Z., 182
 Kourentzes, N., 159
 Koursaros, D., 75, 157
 Kourtellos, A., 142, 167
 Koutra, V., 33
 Kovacs, L., 13
 Kovalchik, S., 75
 Kozubowski, T., 6
 Kral, C., 20
 Krapf, D., 6
 Krasnopjorovs, O., 79
 Kresin, C., 59
 Krivobokova, T., 23, 151
 Krupskiy, P., 85, 86
 Kruse-Becher, R., 82
 Kubo, T., 116
 Kuehn, S., 194
 Kuendig, P., 60
 Kuipers, J., 154
 Kulkarni, A., 52
 Kumavat, M., 130
 Kundu, S., 126
 Kunst, R., 83
 Kurbucz, M., 13
 Kurt, E., 115
 Kurth, J., 109
 Kurtz-Garcia, R., 187
 Kusano, S., 93
 Kutta, T., 118, 158
 Kvaloy, J., 19
 Kyriakou, I., 28
 Kyrilis, I., 92
- Laborda, J., 180
 Lacaza, R., 78
 Laha, N., 97
 Lahiri, P., 186
 Lahiri, S., 57
 Lallemand, D., 155
 Lalli, M., 88
 Lam, C., 136
 Lamarche, C., 30
 Lambardi di San Miniato, M., 130, 131
 Lange, F., 50
 Lange, K., 12
 Langrock, R., 74, 193
 Lanino, L., 68
 Larsen, B., 142
 Latino, C., 120
 Latuszynski, K., 125
 Laumer, S., 115
 Laumont, C., 183
 Laurent, S., 76
 Laureti, T., 110, 144
 Laurini, F., 68
 Lawson, A., 38
 Lazarus, E., 50
- Le, C., 192
 Lechner, M., 143
 Leclercq-Samson, A., 69
 Lee, C., 181
 Lee, D., 98, 107
 Lee, H., 173
 Lee, J., 29, 42, 80, 82, 102, 108, 117, 140
 Lee, K., 26, 62, 139
 Lee, M., 82, 90
 Lee, Q., 30
 Lee, S., 104, 109, 118
 Lee, Y., 56
 Legramanti, S., 127
 Lehmann, S., 81
 Lehoucq, R., 32
 Leinwand, B., 173
 Leisen, F., 155, 184
 Lemster, S., 50
 Leng, C., 156
 Leon, A., 24
 leoni, I., 28
 Lepore, A., 23
 Lesellier, M., 114
 Letixerant, P., 82
 Levantesi, S., 5, 6
 Levin, K., 37, 106, 173, 185
 Levina, L., 58, 125
 Levis, A., 138
 Levrard, C., 152
 Ley, C., 18
 Lhaut, S., 11
 Li, B., 34, 105
 Li, C., 131, 139, 141
 Li, D., 32
 Li, F., 35
 Li, G., 25
 Li, J., 132, 137, 174, 176
 li, L., 189
 Li, M., 136
 Li, Q., 169
 Li, R., 123, 172
 Li, S., 69, 148, 156
 Li, T., 26, 66, 69, 82, 192
 Li, W., 126
 Li, X., 118
 Li, Y., 29, 111, 133, 164
 Li, Z., 34, 69, 174
 Lieber, J., 114
 Lieberman, O., 29
 Liesenfeld, R., 178
 Likitapiwat, T., 147
 Lila, E., 126
 Lim, J., 163
 Lin, E., 67
 Lin, H., 14
 Lin, L., 122, 123
 Lin, T., 138
 Lin, Z., 182
 Lindquist, M., 193
 Linero, A., 164
 Linton, O., 145
 Lipiecki, A., 159
 Lippi, M., 179
 Lisi, F., 165
 Liu, B., 97
 Liu, C., 41, 169
- Liu, F., 4
 Liu, J., 33
 Liu, K., 136
 Liu, P., 32
 Liu, R., 39, 92
 Liu, X., 133
 Liu, Y., 26, 31, 183
 Liu, Z., 35
 Livingstone, S., 90, 174
 Llorens-Terrazas, J., 28
 Llosa, C., 32
 Lloyd, S., 3
 Lo Cascio, I., 12, 63
 Lo, C., 60, 81
 Lo, S., 71
 Loboda, D., 7
 Loh, Y., 51
 Long, Q., 126
 Long, T., 51
 Lonn, R., 85
 Lopes, H., 86
 Lopes, M., 161
 Lopetuso, E., 150, 179
 Lopez Iturriaga, F., 65
 Lopez Lira, A., 101
 Lopez Oriona, A., 103
 Lopez Pintado, S., 103
 Loria, F., 3, 121
 Lorusso, M., 157
 Lourenco, V., 73
 Loyal, J., 126
 Lu, C., 154
 Lu, L., 30
 Lu, Q., 136
 Lu, S., 43
 Lu, W., 152
 Lu, X., 98
 Lu, Y., 22, 44, 123
 Luati, A., 63, 135
 Lubberts, Z., 106, 141, 174
 Lucarno, R., 68
 Lucchetta, M., 120
 Luciani, M., 3
 Luedtke, A., 38
 Luger, R., 60
 Lukemire, J., 142
 Lukens, L., 177
 Lukic, Z., 81
 Lunde, R., 57, 125
 Luo, C., 123
 Luo, R., 103
 Luo, W., 34
 Luo, Y., 111, 124
 Lupi, C., 65, 77
 Lupporelli, M., 154, 163
 Lupu, R., 85
 Lyhagen, J., 164
 Lyrio, M., 49
 Lyu, G., 112
 Lyzinski, V., 106, 141, 174
- Ma, H., 194
 Ma, J., 44, 71
 Ma, L., 9, 141
 Ma, S., 35
 Ma, T., 32
 Ma, Y., 33
- Ma, Z., 103
 Maathuis, M., 174
 Macaulay, V., 189
 MacDonald, P., 173
 Machado, L., 77, 110
 Maciak, M., 81
 MacKinnon, J., 162
 Macri-Demartino, R., 75
 Madari, Z., 13
 Madrid Padilla, C., 104
 Maehashi, K., 91
 Maeng, H., 133
 Maffei Faccioli, N., 100
 Maggi, M., 65
 Maggioni, G., 68
 Magni, G., 5, 99
 Maheshwari, A., 130
 Maignant, E., 152
 Maih, J., 3
 Mailhot, M., 134
 Maillet, R., 21
 Maitra, R., 18, 31
 Majewski, A., 109
 Majumder, R., 176
 Majumder, S., 31
 Majumder, T., 140
 Makarova, S., 85
 Makhdoom, W., 93
 Malevich, N., 123
 Malik, W., 72
 Malinsky, D., 58
 Mallick, B., 122
 Maluf, Y., 76
 Mameli, V., 130, 194
 Mammadli, R., 89
 Mamon, R., 113
 Manafi Neyazi, A., 29
 Maneesoonthorn, W., 47
 Mankad, S., 185
 Manner, H., 21
 Mantoan, G., 3
 Mantziou, A., 12
 Manu, M., 51
 Maraj-Zygmatis, K., 6
 Marbac, M., 107
 Marceau, E., 47, 100
 Marcellino, M., 157
 Marchese, M., 28
 Marchetti, S., 144
 Marcon, G., 153
 Marcoux, M., 113
 Margaritella, L., 134
 Mariani, F., 92
 Marin, J., 69
 Marino, M., 136, 162, 163
 Marion, Z., 89
 Marion, P., 86
 Markatou, M., 70
 Marks-Anglin, A., 123
 Marsman, M., 105, 148
 Martella, F., 136, 162
 Martin, R., 195
 Martin-Utrera, A., 101
 Martinez Martinez, A., 64
 Martinez, B., 159
 Martinez, M., 79
 Martinez-Beneito, M., 39

- Martinez-Garcia, E., 114
 Martini, T., 41
 Martins Bianco, L., 188
 Martos, G., 135
 Maruotti, A., 147
 Masci, C., 14
 Maso, S., 3
 Massacci, D., 44
 Massol, O., 167
 Mastrantonio, G., 53
 Masuda, H., 150
 Matabuena, M., 53
 Matsuda, T., 160, 178
 Matsuda, Y., 60
 Matteo, B., 23
 Mauerer, I., 159
 Mauri, L., 155
 Maxand, S., 132
 May, C., 123
 Mayekawa, S., 108
 Mayo-Iscar, A., 68, 76
 Mayr, A., 130
 Mazo, G., 117
 Mazur, S., 164
 McClean, A., 138
 McDonald, I., 110
 McGee, G., 172
 McGough, O., 185
 McInerney, A., 77
 McIntyre, S., 26, 45
 McKennan, C., 51, 141
 McLachlan, G., 41
 McMillan, P., 177
 McNealis, V., 17
 McNicholas, P., 163
 Medovikov, I., 176
 Meglioli, F., 66
 Mehta, S., 55
 Meilan-Vila, A., 18
 Melly, B., 30
 Melnykov, Y., 166
 Melosi, L., 165
 Mena, R., 64
 Menardi, G., 18
 Mendoza, E., 21
 Meng, W., 138
 Mentch, L., 121
 Mercuri, L., 14, 67, 99
 Merfeld, J., 186
 Mesa Romero, Q., 112
 Messer, K., 138
 Messer, P., 186
 Metulini, R., 87
 Meyer, N., 102
 Meyer, R., 62
 Miao, D., 61
 Michael, S., 166, 167
 Michaelides, D., 171, 182
 Michaelides, M., 100
 Michail, N., 75
 Michetti, E., 65
 Micune, V., 49
 Mies, F., 7
 Miffre, J., 11
 Migliorati, S., 16
 Mikosch, H., 45
 Mikosch, T., 29
 Millard, S., 128, 188
 Miller, C., 72
 Milosevic, B., 81
 Minelli, G., 53
 Minervini, L., 146
 Mingoli, G., 46
 Mira, A., 87
 Miranda Dominguez, O., 142
 Mishra, K., 129
 Mitchell, J., 26, 45, 142
 Mitjans, M., 78
 Mitra, N., 35, 56
 Mittnik, S., 10, 98
 Miyakawa, D., 91
 Mlinar, T., 143
 Modugno, M., 100
 Moffatt, P., 66
 Moghaddam, S., 19
 Moghimbeygi, M., 18
 Mogliani, M., 25
 Mohammadi, R., 105, 153
 Molinari, R., 59
 Momozaki, T., 178
 Mondal, D., 18
 Mondal, P., 131
 Montanes, A., 149
 Montanino, A., 48, 134
 Monteforte, L., 27
 Monteiro, M., 18
 Moodie, E., 17
 Moore, J., 19
 Mora Valencia, A., 24
 Mora, J., 24
 Moradi Rekabdararkolaee, H., 116
 Mordant, G., 1
 Moreira Lara, I., 48
 Moreira, M., 5
 Morelli, A., 149
 Mori, A., 3
 Mori, L., 99
 Morimoto, T., 10
 Morio, J., 60
 Morita, H., 63, 165
 Morita, R., 92
 Moriyama, T., 161
 Mork, D., 17
 Moro, S., 53
 Morota, G., 72
 Morris, C., 24
 Morrone, A., 63
 Morzywolek, P., 38
 Moura, R., 161
 Mousavi, P., 62
 Moutzouris, I., 11, 28
 Mowry, E., 59
 Mozdzen, A., 17
 Mu, X., 25
 Mu, Y., 131
 Mueller, T., 50
 Muhinyuza, S., 164, 181
 Mukherjee, A., 112
 Mukherjee, D., 97
 Mukherjee, R., 54
 Mukhopadhyay, I., 131, 189
 Mukhoti, S., 84, 169
 Muni Toke, I., 101, 150
 Munko, M., 39, 96
 Munoz-Elguezabal, J., 79
 Murach, M., 85
 Murakami, H., 112
 Murphy, K., 90
 Murray Kearney, L., 160
 Murtazashvili, I., 66
 Murua, A., 18
 Mutti, A., 100
 Myers, J., 32
 Mylona, K., 71, 171
 Naboka-Krell, V., 12, 108
 Naef, J., 24
 Nagasaka, S., 160
 Nai Ruscone, M., 16
 Nakagawa, T., 38, 160, 178
 Nakajima, J., 3, 67
 Nakamura, K., 159
 Nalisnick, E., 121
 Nandy, S., 58
 Nane, T., 189
 Nareklishvili, M., 119
 Narita, S., 88
 Nasini, S., 143
 Nasri, B., 86
 Nathoo, F., 183
 Natsiopoulos, K., 65
 Naulet, Z., 188
 Nava, A., 93
 Ndonfack Zango, C., 113
 Neal, M., 163
 Negeri, Z., 144
 Negrete Gallego, R., 70
 Nelson, B., 183
 Neocleous, T., 131
 Nersisyan, L., 49
 Neuhierl, A., 149
 Neuwirth, S., 45
 Newey, W., 5
 Newhouse, D., 186
 Ng, C., 111
 Nguyen, H., 86
 Nguyen, L., 57
 Nguyen, M., 155
 Nguyen, T., 49, 73
 Nguyen, V., 88
 Nicol, F., 129
 Nicolau, J., 83
 Nicolussi, F., 153
 Nielsen, J., 62
 Nielsen, M., 162
 Nieto Delfin, M., 144
 Niguez, T., 24
 Nikolaishvili, G., 115
 Nikoloulopoulos, A., 10
 Ning, N., 192
 Nipoti, B., 188
 Nishi, I., 10
 Nishida, K., 21, 78
 Nishiyama, T., 88
 Nitta, K., 112
 Nixon, G., 195
 Nnanatu, C., 124
 Nodehi, A., 18
 Nolan, J., 102
 Norouzirad, M., 160
 Northrop, P., 120
 Nortier, B., 147
 Novelli, M., 92
 Nuesing, L., 166
 Nunes, M., 15, 74
 Nutarelli, F., 88
 O'Neill, E., 151
 Oates, C., 9
 Odet, A., 107
 Oecal, M., 158
 Oeztuerk, S., 132
 Oganisian, A., 56
 Ogihara, T., 94
 Ogutu, J., 73
 Oh, H., 103, 153
 Ohnishi, Y., 35
 Oishi, R., 63
 Okabe, M., 109
 Okahara, H., 38
 Okano, E., 159
 Oliva, I., 99
 Ombao, H., 103
 Omer, T., 142
 Ommen, D., 166
 Omori, Y., 8
 Orbe, S., 92
 Otava, M., 71
 Otneim, H., 164
 Otranto, E., 47
 Otto, S., 158
 Ouimet, F., 145
 Ovaskainen, O., 181
 Owyang, M., 96, 115, 163
 Oya, A., 79
 Paccagnini, A., 63
 Paci, L., 107
 Pacifici, C., 31
 Packham, N., 113
 Padoan, S., 11, 30, 31, 151
 Paganin, S., 36
 Page, G., 36, 181
 Pak, D., 152
 Palomba, G., 46
 Palumbo, B., 23
 Palumbo, D., 42, 83
 Pan, Y., 128
 Panaretos, V., 152
 Panarotto, A., 39
 Pandolfo, G., 70
 Panorska, A., 6
 Pantelidis, T., 109
 Pantoja, K., 44
 Panzera, A., 18, 117
 Papapostolou, N., 11
 Papastamoulis, P., 16
 Papatsouma, I., 182
 Pappada, R., 134
 Pappalardo, L., 88
 Paraskevopoulos, I., 166
 Parast, L., 138
 Paredes, J., 24
 Pareek, S., 59
 Parikh, H., 190
 Park, C., 38, 127, 143
 Park, J., 107

- Park, S., 74, 153, 156
 Park, Y., 106, 125, 174
 Parker, T., 30
 Parla, F., 12, 63
 Parmigiani, G., 189
 Pascaru, G., 91
 Pasquaretta, C., 107
 Pastor, I., 65
 Pastorino, C., 22
 Patel, N., 165
 Patelli, L., 153
 Pateras, K., 20
 Paterlini, S., 195
 Pathlavath, A., 169
 Paul, J., 82
 Paul, S., 41
 Pauphilet, J., 54
 Pavlidis, E., 114
 Pavliotis, G., 67
 Pavlova, L., 75
 Pecorelli, N., 14
 Pedemonte, M., 142
 Pedroni, P., 183, 184
 Peignon, J., 45
 Peiris, R., 4
 Peiris, S., 63
 Pellegrino, F., 95
 Peng, B., 4
 Peng, P., 76
 Peng, X., 98
 Peng, Y., 123, 172
 Pensky, M., 37
 Peralta Alva, A., 165
 Perchiazzo, A., 14, 99
 Perez, J., 85
 Perez, T., 2
 Peri, I., 22
 Perote, J., 24
 Perotti, R., 165
 Perrakis, K., 16
 Perron, B., 48
 Perry, A., 17
 Perusquia, J., 64
 Peruzzi, A., 42
 Pesaran, H., 1, 5
 Pesavento, E., 27
 Pesta, M., 81, 145
 Pestova, B., 145
 Peters, M., 140
 Peterson, C., 59, 188
 Petrasek, L., 108
 Pettenuzzo, D., 179
 Pfeil, T., 156
 Phylaktis, K., 11, 12
 Piccarreta, R., 36
 Pick, A., 5
 Piepho, H., 72
 Pievatolo, A., 69
 Pigeon, M., 100
 Piger, J., 96, 163
 Pignata, S., 89
 Pigoli, D., 143, 158
 Pini, A., 52
 Pinocchio, N., 68
 Pintado, M., 86
 Piperigou, V., 129
 Pipis, C., 37
 Piscopo, G., 5, 6
 Pitt, M., 62
 Pittavino, M., 15
 Pizarro-Irizar, C., 159
 Poggi, J., 61
 Poggio, D., 52, 53
 Politis, D., 23
 Poncela, P., 135
 Pons, M., 30
 Poonia, N., 129
 Poplawski Ribeiro, M., 180
 Porqueddu, M., 63
 Portier, B., 61
 Portugal, P., 83
 Potjagailo, G., 177
 Pouliasis, P., 11
 Pouliot, W., 114
 Pozuelo Campos, S., 70
 Pozza, F., 9
 Prates, M., 132, 137
 Priebe, C., 58, 74, 106, 125, 174
 Prifti, O., 134
 Primiceri, G., 100
 Proietti, T., 63, 95
 Prokhorov, A., 50
 Prosdocimi, I., 120
 Prostmaier, B., 12
 Prueser, J., 48
 Pua, A., 132
 Puc, A., 7
 Pupo West, T., 194
 Purkayastha, S., 14
 Qi, T., 141
 Qi, Z., 42
 Qiao, W., 105
 Qiao, X., 126
 Qin, L., 44
 Qin, X., 169
 Qiu, G., 141
 Quach, N., 138
 Quaini, A., 148
 Queiroz, F., 76
 Quetti, F., 15
 Quintana, F., 36, 137
 Quiroz, M., 47
 Quiroz, Z., 132
 Rabitti, G., 153
 Rabonza, M., 155
 Raffinetti, E., 87
 Raftapostolos, A., 26
 Raggi, D., 84
 Rahbek, A., 29
 Rahman, N., 72
 Rainone, E., 183
 Ramdas, A., 189
 Ramezan, R., 138
 Ramirez-Silva, I., 64
 Ramos, A., 149
 Ramsay, K., 103
 Ranaldo, A., 150
 Rancoita, P., 14
 Randolph, A., 142
 Rao, J., 16
 Raponi, V., 101
 Raposo, P., 83
 Rastelli, R., 154
 Ravazzolo, F., 83
 Ray, S., 72
 Rayan, S., 121
 Reale, M., 21
 Reeve, H., 190
 Reich, B., 90, 176
 Reinert, G., 37
 Reis, H., 83
 Reluga, K., 186
 Remillard, B., 86, 143
 Ren, S., 74
 Rennspies, J., 29
 Resin, J., 168
 Rezitis, A., 66
 Riabiz, M., 9
 Ricardo, I., 179
 Ricciotti, L., 147
 Ricco, G., 27
 Richard, F., 135
 Rieger, O., 165
 Rigon, T., 56, 188
 Rios, N., 104
 Risso, D., 20
 Risstad, M., 157
 Ritz, A., 129
 Rivera, N., 39
 Rivieccio, G., 89
 Rizzelli, S., 11, 30, 31, 151
 Roberts, G., 125, 150, 175
 Robertson, D., 26, 55
 Robles, M., 31
 Roca Pardinias, J., 77
 Rochanahastin, N., 147
 Rodrigues, P., 83
 Rodriguez, C., 157
 Rodriguez-Caballero, V., 111
 Rodriguez-Poo, J., 2, 3, 168
 Roettger, F., 120
 Rogantini Picco, A., 165
 Rohe, K., 113
 Romagnoli, L., 77
 Romanak, M., 145
 Romano, G., 154
 Romero Bejar, J., 112
 Rosa, P., 106
 Rosci, E., 72
 Rose, C., 183
 Rossell, D., 140, 153
 Rossi, E., 46, 149, 150
 Rossi, F., 29
 Rossi, L., 27
 Rossini, L., 86
 Rothman, A., 82
 Roviello, A., 5
 Rowlinson, E., 34
 Roy, A., 59, 97, 169
 Roy, S., 74, 84, 122, 140, 184
 Rroji, E., 14, 99
 Ruan, F., 25
 Rubaszek, M., 85
 Ruberu, T., 189
 Rubin-Delanchy, P., 113
 Rubini, L., 185
 Rudolph, K., 17, 190
 Ruecker, M., 145
 Ruggiero, M., 38
 Ruiz, E., 1, 111
 Ruiz-Castro, J., 77
 Russell, T., 114
 Russo, A., 68, 147
 Russo, M., 41, 124
 Rust, G., 175
 Saadaoui, J., 146
 Sabbatucci, R., 179
 Sabbioni, E., 15
 Sadr, N., 134
 Saefken, B., 129
 Safari-Katesari, H., 187
 Safikhani, A., 119
 Safo, S., 191
 Saha, E., 128
 Saha, S., 59
 Sahay, A., 148
 Sahu, S., 124
 Saine, M., 123
 Sakhanenko, L., 144
 Sakiyama, T., 4
 Sala, L., 165
 Salmaso, F., 52
 Salvagnin, C., 87
 Salvati, N., 186
 Salvini, N., 110, 144
 Samadi, S., 116, 187
 Samarawickrama, N., 51
 Samartsidis, P., 15
 Sancetta, A., 26
 Sanchez-Betancourt, L., 90
 Sanchez-Sellero, C., 161
 Sanchez-Torres, J., 79
 Sanchis-Marco, L., 31
 Sandberg, R., 157, 181
 Sanfelici, S., 47
 Sangalli, L., 23, 69
 Sangiovanni, G., 53
 Sankaran, K., 44
 Sanna Passino, F., 12, 58
 Sant, J., 9
 Santoro, A., 68
 Santos, I., 115
 Santucci de Magistris, P., 150
 Sarabia, J., 2
 Sarafidis, V., 45
 Sarault, C., 135
 Sarfo Fosu, E., 191
 Sarhadi, A., 72
 Sariev, H., 87
 Sarisoy, C., 61
 Sarkar, S., 117
 Sarrocco, G., 120
 Sasaki, T., 91
 Sato, Y., 112
 Satten, G., 59
 Satterthwaite, T., 142
 Saunders, C., 166
 Sauri, O., 150
 Sauta, E., 68
 Savevski, V., 68
 Savva, C., 75
 Sawaya, K., 8
 Sayer, T., 113
 Scaccia, L., 65

- Scaffidi Domianello, L., 47
 Scandurra, G., 146
 Schafer, T., 136
 Schaffer, M., 82
 Schauburger, G., 130
 Scheepers, B., 176
 Scherpenhuijzen, N., 194
 Scheufele, R., 83
 Schiavon, L., 36
 Schiavone, F., 89
 Schildcrout, J., 142
 Schisa, V., 118
 Schmid, T., 186
 Schnaitmann, J., 76
 Schnattinger, P., 177
 Schnurbus, J., 132
 Schnurr, A., 8
 Scholze, F., 7
 Schoonhoven, M., 105, 153
 Schuessler, R., 81
 Schuettler, J., 108
 Schulz, C., 121, 168
 Schumann, M., 109
 Schwabe, R., 123
 Schweinberger, M., 37
 Schwob, M., 137
 Scognamiglio, S., 89
 Scott, M., 72
 Secondi, L., 144
 Seewald, N., 56
 Segala, C., 102
 Segers, J., 11
 Seidlitz, J., 142
 Sekine, T., 63
 Sekulovski, N., 148
 Seleznova, M., 86
 Selland Kleppe, T., 60
 Semeraro, P., 100
 Semmler, W., 98, 99
 Sen, A., 186
 Sen, R., 103
 Senn, E., 194
 Sentana, E., 115
 Senturk, D., 122
 Seo, M., 5
 Serra, F., 6
 Serrubeco Marques, E., 167
 Seshagiri, A., 184
 Sestelo, M., 77, 110
 Seto, H., 108
 Severino, F., 135
 Shah, M., 93
 Shakhnov, K., 12
 Shaman, J., 177
 Shang, L., 191
 Sharifi Kiasari, M., 106
 Sharifvaghefi, M., 5
 Sharma, P., 84
 Shelley, M., 74
 Shen, C., 125, 173
 Shen, H., 173
 Shen, S., 174
 Shen, X., 41
 Shen, Z., 9
 Sheng, A., 191
 Sheng, S., 82
 Shestopaloff, A., 86, 90
 Shi, C., 42, 82
 Shi, J., 111
 Shi, Y., 11, 28
 Shields, K., 26
 Shimadzu, H., 78
 Shimizu, K., 162
 Shimokawa, A., 79
 Shin, H., 170
 Shin, Y., 64
 Shioji, E., 64
 Shiotani, T., 67
 Shiroff, T., 142
 Shokoohi, F., 91
 Short, N., 172
 Shrimali, G., 120
 Shukla, A., 192
 Shulman, D., 177
 Shushi, T., 91
 Si, Y., 132
 Sibillo, M., 5, 99
 Sieg, F., 96
 Signorelli, M., 146
 Sigrist, F., 60, 78, 93
 Silbernagel, A., 19
 Silvennoinen, A., 98
 Simone, R., 61
 Simoni, A., 25, 45, 95
 Simons, J., 145
 Simpson, A., 166, 167
 Simpson, E., 120
 Sina, A., 120
 Singh, A., 124
 Sinha, D., 122, 175
 Sinha, S., 33
 Sisson, S., 11
 Skhosana, S., 128
 Skowronek, K., 6
 Skrobotov, A., 50
 Smadu, A., 99
 Smedts, K., 49
 Smeekes, S., 88
 Smith, L., 66
 Smith, M., 47
 Smith, R., 1, 5, 145
 Smyth, R., 146
 So, M., 3
 Soale, A., 133
 Soberon, A., 3, 168
 Sochaniwsky, A., 163
 Soegner, L., 43
 Soehl, J., 8
 Solea, E., 34
 Sonabend, A., 97
 Sone, T., 91
 Song, D., 118
 Song, H., 125
 Song, J., 191
 Sonobe, N., 178
 Soques, D., 163
 Sorel, A., 74, 107
 Sornnil, N., 147
 Soto, C., 52
 Sottosanti, A., 20
 Spano, D., 9
 Speagle, J., 32
 Sperlich, S., 21
 Spicker, D., 103
 Spieker, A., 54
 Spini, P., 119
 Sprincenatu, M., 10
 Sriram, K., 84
 Stafoggia, M., 72
 Stanton, R., 50
 Staszewska-Bystrova, A., 12, 108
 Steel, M., 43, 162
 Stefanelli, K., 6
 Stefani, I., 14, 99
 Steiner, G., 162
 Steland, A., 7, 155
 Stenning, D., 32
 Stensrud, M., 127
 Stewart, J., 141
 Steyn, A., 145
 Stingo, F., 124
 Stival, M., 36
 Stoecker, A., 23, 152
 Stoeve, B., 164
 Stoner, O., 131
 Stove, B., 146
 Strawderman, R., 17
 Strenger-Galvis, D., 159
 Striaukas, J., 24, 45
 Stuart, E., 190
 Stupfler, G., 10, 151
 Su, H., 171
 Sucarrat, G., 95
 Sudderth, E., 140
 Sugasawa, S., 8, 38
 Sulem, D., 153
 Sun, X., 35
 Sun, Y., 103, 123, 139
 Susmann, H., 17
 Sussman, D., 174
 Suzuki, S., 111
 Svetlosak, A., 108
 Sykulski, A., 13
 Szafranek, K., 85
 Szerszen, P., 185
 Szymanski, G., 21
 Szymczak, S., 50
 Tahata, K., 160
 Takabatake, T., 7
 Takahashi, H., 91
 Takahashi, K., 20, 78
 Tambalotti, A., 122
 Tanda, A., 88
 Tang, J., 26
 Tang, S., 185
 Tang, T., 97
 Tang, W., 192
 Tang, X., 42, 132, 139, 191, 192
 Tang, Y., 34, 62
 Tao, R., 142
 Tarantola, C., 63
 Taschini, L., 68
 Tavakoli, S., 23
 Tavin, B., 176
 Taylor, R., 142, 162
 Tchetgen Tchetgen, E., 127
 Tebaldi, P., 114
 Tegge, A., 170
 Tejeria-Martinez, M., 2
 Telcs, A., 13
 Telg, S., 46
 Tentori, C., 68
 Terada, Y., 111
 Terasvirta, T., 98
 Tervo-Clemmens, B., 142
 Testa, L., 52
 Tetereva, A., 85, 151
 Thiede, P., 12
 Thiery, A., 175
 Thoeni, P., 68
 Thomas, A., 45, 46, 167
 Thompson, J., 70, 110
 Thongtai, C., 147
 Thorsen, I., 146
 Thulliez, E., 61
 Tian, Y., 141
 Timmermann, A., 179
 Todisco, G., 68
 Tomarchio, S., 47
 Tomilina, E., 117
 Tommasi, C., 123
 Tonaki, Y., 93, 102
 Tong, G., 20
 Tong, J., 123
 Tong, M., 3
 Tong, T., 122
 Tong, X., 119
 Toninelli, D., 153
 Topaloglou, N., 180
 Toparkus, A., 96
 Torelli, N., 75
 Torgovitsky, A., 114
 Tran, J., 19, 60
 Tran, M., 4
 Trapani, L., 116, 133, 179
 Trapote-Reglero, L., 76
 Trede, M., 81
 Tremayne, A., 111
 Tria, F., 88
 Triantafyllou, A., 180
 Trinca, L., 71
 Trinh, D., 177
 Trosset, M., 58
 Trovato, G., 185
 Tsagris, M., 133
 Tsai, H., 67
 Tsujimoto, T., 54
 Tsukahara, H., 10
 Tsyawo, E., 66, 133
 Tu, X., 138
 Tudorascu, D., 142
 Tutz, G., 159
 Tuvaandorj, P., 48
 Tzivanakis, N., 13
 Uberti, P., 22, 156
 Ubezio, M., 68
 Uchida, M., 93, 101
 Uematsu, Y., 104
 Ullmann, T., 50
 Umlandt, D., 151, 156
 Uniejewski, B., 194
 Uppal, R., 101
 Usseglio-Carleve, A., 151
 Vaida, F., 138

- Vaillancourt, J., 143
 Valade, T., 8
 Valente, E., 6
 Vallejos, R., 30
 Vallevik, V., 181
 Valvason, C., 21
 van Arem, K., 8
 Van Binsbergen, J., 101
 van de Werve, I., 194
 van den Bergh, D., 148
 van den Boogaart, K., 2
 Van Den Broeke, M., 143
 van der Laan, M., 186
 van der Molen Moris, J., 57
 van der Oord, J., 46
 van Eeuwijk, F., 73
 van Heerwaarden, J., 154
 Van Keilegom, I., 96
 van Kesteren, E., 59
 van Spronsen, J., 63
 van Wieringen, W., 105
 Vandekar, S., 142
 Vanduffel, S., 48
 Vannucci, M., 105, 127, 140
 Vantini, S., 52, 70
 Varelas, S., 92
 Varin, C., 117
 Vasdekis, G., 9
 Vats, A., 130
 Vats, D., 192
 Vaudree, L., 46
 Veiga, H., 69
 Velasco, S., 75
 Velasquez-Gaviria, D., 95
 Veldhuis, S., 28
 Ventouri, A., 45
 Ventura, D., 53
 Veraart, A., 55
 Veraart, L., 113
 Verster, T., 176
 Vespignani, J., 146
 Vidoni, P., 131, 194
 Viechnicki, P., 141
 Vienken, M., 74, 193
 Vilandt, F., 29
 Villa, C., 64, 184
 Villani, M., 62, 129
 Villanueva, N., 77, 110
 Villar, S., 55
 Villejo, S., 72
 Violante, F., 95
 Virbickaite, A., 86
 Viroli, C., 16
 Visagie, J., 145
 Vitale, D., 66
 Vo, T., 17
 Vogels, L., 105, 153
 Vogt, M., 145
 Voirol, L., 59
 von Tycowicz, C., 152
 Vriz, G., 68
 Waagepetersen, R., 170
 Wade, G., 98
 Wade, S., 36
 Wagner, K., 132
 Wagner, M., 28, 43
 Wahed, A., 152
 Waite, T., 34
 Walden, J., 49
 Waldorp, L., 148
 Wallace, M., 121
 Walsh, C., 145
 Walwyn, R., 124
 Wang, B., 126
 Wang, C., 4, 77, 158
 Wang, G., 193
 Wang, H., 79, 97, 170, 171
 Wang, J., 24–26, 42, 169, 192
 Wang, L., 97
 Wang, O., 123
 Wang, Q., 42, 186
 Wang, R., 53, 152
 Wang, S., 158, 193
 Wang, T., 3, 167
 Wang, X., 136, 171
 Wang, Y., 23, 51, 106, 139
 Wang, Z., 47
 Ward, O., 32
 Wasserman, L., 138
 Watanabe, T., 67
 Weale, M., 162
 Weaver, C., 128
 Webb, A., 71
 Webb, M., 162
 Weber, M., 141
 Weerakoon, K., 186
 Wegener, C., 82
 Wei, C., 66
 Wei, Y., 22, 57
 Weidner, M., 44
 Weiss, C., 19
 Weissbach, R., 96
 Wen, L., 172
 Wen, Q., 128
 Weron, R., 159, 194
 Wertz, T., 17
 West, M., 3
 White, A., 163
 Wied, D., 158, 178
 Wieladek, T., 162
 Wikle, C., 81
 Wikle, N., 187
 Wilkie, C., 72
 Williams, N., 17
 Williamson, D., 147
 Wilms, I., 88
 Wilson, J., 123
 Wilson, K., 184
 Wilson, L., 54
 Winkelmann, D., 74, 193
 Winker, P., 12, 108
 Wiper, M., 72
 Witkovsky, V., 81
 Witten, D., 185
 Witulska, J., 6
 Wixson, T., 137
 Wolf, E., 49
 Wolf, M., 43
 Wolfram, D., 168
 Wolny-Dominiak, A., 61
 Wortham, C., 81
 Wozniak, T., 121
 Wright, S., 26
 Wrobel, J., 142
 Wu, C., 139
 Wu, D., 124
 Wu, P., 45
 Wu, Q., 128
 Wu, W., 108, 156
 Wu, X., 105
 Wu, Y., 128
 Wylomanska, A., 6, 7
 Xi, M., 9
 Xiao, H., 137
 Xiao, L., 128
 Xie, C., 172
 Xie, F., 124
 Xie, Y., 158
 Xiong, J., 142
 Xiong, Y., 125, 139
 Xu, G., 170
 Xu, H., 103, 123
 Xu, J., 140
 Xu, M., 132
 Xu, Q., 88
 Xu, S., 170
 Xu, W., 37
 Xu, X., 33
 Xu, Y., 9, 76
 Xue, F., 35
 Xue, H., 139
 Xue, L., 70, 105
 Yaari, R., 177
 Yadohisa, H., 109
 Yamagata, T., 115
 Yamaguchi, H., 112
 Yamamoto, M., 78, 108, 109
 Yamanaka, S., 22
 Yamauchi, Y., 8
 Yan, J., 137
 Yan, Y., 4
 Yang, C., 5
 Yang, H., 66, 144
 Yang, J., 67
 Yang, S., 79
 Yang, T., 67
 Yang, X., 152
 Yang, Y., 123, 180
 Yano, T., 8
 Yao, A., 129
 Yao, F., 167
 Yao, Q., 80
 Yao, W., 188
 Yata, K., 7
 Yau, C., 111
 Ye, C., 90
 Ye, S., 29
 Ye, Z., 34
 Yeganegi, M., 22
 Yeon, H., 103
 Yi, M., 53
 Yoshida, J., 101
 Yoshida, N., 101, 150
 Yu, C., 111
 Yu, D., 80
 Yu, J., 48
 Yu, L., 183
 Yu, M., 127
 Yu, S., 106, 139
 Yu, W., 73, 111
 Yuan, Y., 25
 Yun, H., 152
 Yurko, R., 107
 Zaborowski, P., 159
 Zadlo, T., 61
 Zaffaroni, P., 101
 Zahra, N., 93
 Zakiyeva, N., 91
 Zakoian, J., 46, 85, 95
 Zanella, G., 9
 Zanetti, F., 165
 Zanutti, G., 153
 Zapf, R., 157
 Zaroudi, S., 187
 Zarraga, A., 159
 Zazzetti, E., 68
 Zeleneev, A., 119
 Zeng, J., 129
 Zeni, G., 70
 Zetterstrom, S., 55
 Zhang, C., 101, 139
 Zhang, J., 60, 80
 Zhang, L., 69, 76, 136, 190, 191
 Zhang, P., 193
 Zhang, Q., 88
 Zhang, W., 179
 Zhang, X., 104, 138
 Zhang, Y., 26, 35, 41, 57, 82, 117, 173, 176, 177
 Zhang, Z., 11, 12, 111, 172
 Zhao, B., 35
 Zhao, N., 59
 Zhao, Y., 26, 172
 Zheng, C., 64, 114
 Zheng, J., 104
 Zheng, L., 186
 Zhong, C., 29, 80
 Zhong, P., 11
 Zhou, D., 118, 174
 Zhou, H., 42, 139
 Zhou, R., 123
 Zhou, S., 192
 Zhou, W., 26
 Zhou, X., 171
 Zhu, D., 144
 Zhu, H., 71
 Zhu, J., 125
 Zhu, L., 90
 Zhu, M., 60
 Zhu, R., 90
 Zhu, X., 117
 Zhu, Y., 58
 Zhu, Z., 93, 192
 Ziegel, J., 168
 Ziel, F., 168
 Zipunnikov, V., 59
 Zitovsky, J., 54
 Zoe Ricci, F., 140
 Zoli, M., 62
 Zou, F., 39
 Zou, G., 55
 Zou, Y., 54, 166
 Zuo, X., 123

