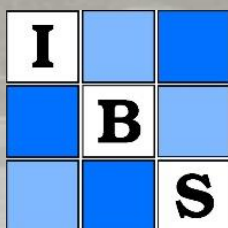


**12th International Conference of the
International Biometric Society's
Eastern Mediterranean Region**

ABSTRACT BOOK





Message from the Chair of the Local Organizing Committee

Dear Colleagues,

On behalf of the organizing committee, I warmly invite you to the 12th International Conference of the International Biometric Society's Eastern Mediterranean Region (EMR 2023) event that will be held in İzmir, the beautiful city of Türkiye, on May 8-11, 2023.

We are happy to be able to hold this event face-to-face after a long time. We are trying to prepare an event that is rich in both scientific and social aspects for you. EMR 2023 will include pre-conference tutorials, invited sessions, round tables as well as oral and poster presentations. The event will also include an honorary mini-symposium dedicated to marking the retirement of Professor Refik Burgut and Professor Ergun Karaağaoğlu. Both professors gave huge efforts of promoting Biostatistics in both Türkiye and the EMR region. This symposium will take place in the afternoon of the first day of the event and will include invited talks.

İzmir is the third largest metropolitan city in western Türkiye, as well as a port city on the tentative list of UNESCO World Heritage. With its 8500-year history, İzmir has hosted many civilizations in its geography. Tepekule, which is the oldest and most historical city in the history of the West, was discovered in İzmir. The Temple of Artemis, one of the 7 wonders of the world, is in Ephesus. The conference venue of EMR 2023 is in the center of İzmir, within walking distance of İzmir Bay and the coastline. As part of the social program of the event, a social trip will be organized to Ephesus, The House of the Virgin Mary, and Sirince Village.

We will be very happy and honored to see you among us.

Sincerely,

Gökmen Zararsız
Conference Chair



About the Conference

The International Biometric Society (IBS) is an international society for the advancement of biological science through the development of quantitative theories and the application, development, and dissemination of effective mathematical and statistical techniques. Areas of application of these methods include, for example, agriculture, medicine, bioinformatics, and ecology.

The society welcomes members from biology, mathematics, statistics, and similar fields. The IBS is the principal international body of biostatisticians and biometricians and comprises Regions throughout the world. The EMR is the region that was formed in 2001 at the inaugural conference held in Athens, Greece. The Region covers Cyprus, Egypt, Greece, Israel, Jordan, the Palestinian Authority, Saudi Arabia, and Türkiye.

The latest conference of the EMR was virtual due to the Covid-19 pandemic. The forthcoming conference, the 12th Conference of the Eastern Mediterranean Region (EMR) of the International Biometric Society (IBS) will convene in Izmir, Turkey, on May 8-11, 2023. In addition to the main conference, a mini-symposium will be held before the main conference, which is devoted to the retirement of former EMR presidents **Prof. H. Refik Burgut** (2001-2003) and **Prof. A. Ergun Karaagaoglu** (2010-2011).

The conference will bring together experts in various disciplines from our Region and all over the world and will provide the opportunity to forge new collaborations, discuss issues of common interest, and plan future developments in biometry, biostatistics, statistics, etc. Besides being of great scientific value in itself, scientific collaboration is one of the ways of enhancing meaningful ties between individuals in different countries and promoting positive and peaceful relationships. We look forward to welcoming you in person in the **EMR2023** conference that is as exciting as our previous conferences held in other EMR region countries.



Scientific Committee

Ahmet Öztürk	Erciyes University, Türkiye
Atilla H. Elhan	Ankara University, Türkiye
Bahar Taşdelen	Mersin University, Türkiye
Benjamin Reiser	University of Haifa, Israel
Bella Vakulenko-Lagun	University of Haifa, Israel
Birol Emir	Pfizer Inc & Adjunct Professor at Columbia University, USA
Cemil Çolak	İnönü University, Türkiye
Cengiz Bal	Eskişehir Osmangazi Üniversitesi, Türkiye
Christos T. Nakas	University of Thessaly, Greece
Constantine Gatsonis	Brown University, USA
David Zucker	Hebrew University of Jerusalem, Israel
Dimitris Karlis	Athens University of Economics and Business, Greece
Erdem Karabulut	Hacettepe University, Türkiye
Ferhan Elmali	Katip Çelebi University, Türkiye
Fikret Er	Anadolu University, Türkiye
Geert Molenberghs	Hasselt University, Belgium
Havi Murad	Getner Institute, Israel
Hernando Ombao	King Abdullah University, Saudi Arabia
Itai Dattner	University of Haifa, Israel
İlker Ercan	Uludağ University, Türkiye
Jaroslav Harezlak	Indiana University, USA
Kenan Köse	Ankara University, Türkiye
Konstantinos Fokianos	Cyprus University, Cyprus
Mehmet Orman	Ege University, Türkiye
Meriç Yavuz Çolak	Başkent University, Türkiye
Mithat Gönen	Memorial Sloan Kettering Cancer Center, USA
Ori Davidov	University of Haifa, Israel
Özlem İlk	Middle East Technical University, Türkiye
Philip Reiss	University of Haifa, Israel
Pınar Özdemir	Hacettepe University, Türkiye
Ruth Heller	University of Haifa, Israel
Siddik Keskin	Yüzüncü Yıl University, Türkiye
Urania Dafni	University of Athens, Greece
Ünal Erkokmaz	Sakarya University, Türkiye
Vildan Sümbüloğlu	Sanko University, Türkiye
Yavuz Sanisoğlu	Yıldırım Beyazıt University, Türkiye
Yoav Benjamini	University of Tel Aviv, Israel



Local Organizing Committee

Honorary Chairs

A. Ergun Karaağaoğlu

Hacettepe University, Türkiye (Emeritus)

H. Refik Burgut

Çukurova University, Türkiye (Emeritus)

Chair

Gökmen Zararsız

Erciyes University, Türkiye

Secretary

Necla Koçhan

IBG (Izmir Biomedicine and Genome Center), Türkiye

Organizing Committee

Dinçer Göksülük

Erciyes University, Türkiye

Erdal Coşgun

Microsoft Genomics, USA

Göknur Giner

Walter and Eliza Hall Institute of Medical Research, Australia

Gözde Ertürk Zararsız

Erciyes University, Türkiye

İlker Ünal

Çukurova University, Türkiye

Konstantinos Fokianos

University of Cyprus, Cyprus

Merve Kaşıkçı

Hacettepe University, Türkiye

Mustafa Çavuş

Eskişehir Technical University, Türkiye

Osman Dağ

Hacettepe University, Türkiye

Sevilay Karahan

Hacettepe University, Türkiye

Vilda Purutçuoğlu

Middle East Technical University, Türkiye



Invited Speakers

Keynote Speaker

Vincent Carey Harvard Medical School, USA

Marvin Zelen Memorial Lecture Speaker

Guadalupe Gómez Melis Universitat Politècnica de Catalunya, Spain

Plenary Speakers

Anne-Laure Boulesteix Ludwig Maximilian University of Munich, Germany
Gordon Smyth Walter and Eliza Hall Institute of Medical Research, Australia
Karen Kafadar University of Virginia, USA
Ping Hu National Cancer Institute, USA
Yoav Benjamini Tel Aviv University, Israel

Invited Speakers

Almond Stöcker École Polytechnic Fédéral de Lausanne, Switzerland
Aris Perperoglou GSK, UK
Çağdaş Hakan Aladağ Hacettepe University, Türkiye
David M. Steinberg Tel Aviv University, Israel
Emrah Gecili University of Cincinnati, USA
Geert Molenberghs Catholic University of Leuven, Belgium
Giorgos Bakoyannis Athens University, Greece
Göknur Giner Walter and Eliza Hall Institute of Medical Research, Australia
Havi Murad Gertner Institute, Israel
KyungMann Kim University of Wisconsin-Madison, USA
Malgorzata Bogdan Lund University, Sweden & University of Wrocław, Poland
Mehmet Gönen Koç University, Türkiye
Mehmet Koçak Istanbul Medipol University, Türkiye
Philip Tzvi Reiss University of Haifa, Israel
R. Todd Ogden Columbia University, USA
Ramyar Molania Walter and Eliza Hall Institute of Medical Research, Australia
Sipan Aslan King Abdullah University of Science and Technology, Saudi Arabia
Victor Kipnis National Cancer Institute, USA

Symposium Speakers

H. Refik Burgut Çukurova University, Türkiye (Emeritus)
A. Ergun Karaağaoğlu Hacettepe University, Türkiye (Emeritus)
Benjamin Reiser University of Haifa, Israel
Hamparsum Bozdogan University of Tennessee, USA
Özlem İlk Middle East Technical University, Türkiye
Recai Yücel Temple University, USA
Yahya Laleli Duzen Laboratories Group, Türkiye



Courses

The courses are held on the first day (May 8, 2023) of the EMR Conference. The available courses are listed below:

Course 1

"Introduction to explainable machine learning with examples in healthcare" by Przemyslaw Biecek from Warsaw University of Technology

The aim of the workshop is to introduce participants to explainable artificial intelligence (XAI) methods that can be used to build predictive models and extract knowledge from predictive models. The workshop will combine discussion of the theoretical basis together with examples with code for your own execution. We will use real-world data for a mortality prediction problem for covid or classification problem for heart disease.

The discussed methods are available in many programming languages and various libraries, but the workshop will be based on examples in R using the DALEX library. The scope of the workshop coincides with that of the book Explanatory Model Analysis <https://ema.drwhy.ai/>.

The first part of the workshop is dedicated to exploratory data analysis tools and preparing for modelling. The second part of the workshop is focused on tools for developing predictive models. For the purposes of the example, we will discuss decision trees, random forests and techniques for automatic tuning of random forests. The third part will focus on local model explanation techniques. We will discuss SHAP (Shapley values), break-down and LIME, the most popular methods for local exploration of models. The fourth part will be devoted to global model explanation techniques. We will discuss the permutation importance technique for variables and the Partial Dependence technique. The workshop will be based on material from <https://github.com/BetaAndBit/RML>

Why?

Complex machine learning models are frequently used in predictive modeling. There are a lot of examples for random forest-like or boosting-like models in medicine, finance, agriculture, etc. But who trusts in black boxes? In this workshop we will show why and how one would analyse the structure of the black-box model. This will be a hands-on workshop. In each part there will be a short lecture and then time for practice and discussion. Using the example of analysing a specific dataset, we will show the basics of modelling with tree models. We will then show how to evaluate and analyse such models using XAI techniques. From the packages, we will learn about randomForest, party, mlr3, DALEX, modelStudio and arenar.



Course 2

"Removing unwanted variation from large-scale RNA sequencing data with PRPS" Ramyar Molania and Marie Trussart from Walter and Eliza Hall Institute

Large scale datasets generated by different omics technologies present unique challenges in terms of normalization and integration. This course focuses on expanding biostatistical and bioinformatics methods for such challenges. We will be focusing on the RUV normalization methods, which have shown great promise in dealing with the challenges presented by large scale datasets from TCGA. RUV-PRPS which is a novel strategy (Molania et al, 2023, Nat. Biotech, <https://www.nature.com/articles/s41587-022-01440-w#code-availability>) uses pseudo-replicates of pseudo-samples (PRPS) to normalize RNA-seq data in situations when technical replicate is not available. In this course we will be presenting the new RUV-PRPS package we have been developing, which is a user-friendly R package that enable researchers to run RUV-PRPS method and to visualize diagnostic plots before and after normalization to assess the quality and consistency of their data.

Session 1: Introduction to large-scale RNA sequencing and RUV methods - Theoretical session

- Introduction on Removing Unwanted Variation (RUV) methods and model
- Pseudo-replicates and pseudo-samples approach (Ramyar et al, Nature Biotech, 2023, <https://www.nature.com/articles/s41587-022-01440-w#code-availability>)

Session 2: Identification of unwanted variation in RNA-seq data - Hands on session

- RNA-seq from the Cancer Genome Atlas (TCGA) and their provided normalizations
- RUV-PRPS package with statistical methods to identify unwanted variation:
 - Functions to identify variation in categorical variables: PCA, silhouette coefficient, ARI, ANOVA, vector correlation.
 - Functions to identify variation in continuous variables: Linear regression, correlation.

Session 3: How to apply RUV-PRPS - Hands on session

- Selection of negative control genes
- Function from RUV-PRPS package to create pseudo-replicates of pseudo-sample (PRPS) to correct for library size, batch effects and tumour purity.
- Function from RUV-PRPS package to run RUV-PRPS method.

Session 4: Normalization performance assessment - Hands on session

- RUV-PRPS package with statistical methods to assess the performance of normalization method:
 - Functions to identify variation in categorical variables: PCA, silhouette coefficient, ARI, ANOVA, vector correlation.
 - Functions to identify variation in continuous variables: Linear regression, correlation.
- How unwanted variation can influence down-stream analysis including gene-gene correlation, survival analysis



Course 3

"How to use cloud technologies for reliable and responsible Data Science projects?"

by Erdal Coşgun from Microsoft Genomics, Vincent Carey from Harvard Medical School, and Deniz Ilhan Topcu from Izmir Tepecik Education and Research Hospital

Researchers are using cloud environments for biomedical data sharing, run analysis tools, and collaborate. In this hands-on course, we will cover the following topics:

- Create a reliable and secure cloud environment. Data Sharing and Auto scale of your compute solutions – 60 mins
- Deploy and use Jupyter Lab, VS Code Server on terra.bio – 30 mins
- Selected topics in genomic visualization and analysis with Bioconductor on cloud - 75 mins
- Responsible Data Science use-cases – 15 mins

Requirements:

- Mid-level (200) Linux OS experience
- Mid-Level (200) R programming experience
- Mid-Level (300) Python programming experience
- Experienced in 'Jupyter Notebook/Lab OR Hub' usage with different kernel types (R, Python, Julia, Spark etc.)
- Virtual Machines will be provided in the course, but participants SHOULD BRING THEIR LAPTOPS OR PC.



Roundtable

Roundtable on “The role of statistical experts during the COVID-19 pandemic”

Session organizers and chairs: David Steinberg (Tel Aviv University, Israel) and Geert Molenberghs (Unversiteit Hasselt and KU Leuven, Belgium).

Panelists:

- Arne Bathke (Universität Salzburg, Austria)
- Ralph Brinks (Witten University, Germany)
- Amit Huppert (Gertner Institute and Tel Aviv University, Israel)
- Filomena Maggino (Sapienza Università di Roma, Italy)
- Bhramar Mukherjee (University of Michigan, Ann Arbor, Michigan, USA)

Description

Around the globe, the COVID-19 induced pandemic has placed scientists in policy advisory roles and brought them to the public forum in ways and of an intensity that they are typically not used to. In some countries, but not everywhere, this also involved biostatisticians and epidemiologists. We will have a fine conversation with some of them, organized around the following questions:

- Looking back at efforts in your country to develop evidence-based policies for COVID-19 response, what was the most important contribution of the statistical community?
- What was the most serious missed opportunity?
- Were statisticians a part of the scientific dialog providing information and advice to decision makers and to the public at large?
- Has the statistical community drawn lessons from the experience as to what was, and was not, successful?
- What actions are essential to improve the ability of the statistical community to be part of the decision-making process in the next pandemic, or similar catastrophe, in other words, what are the conclusions towards pandemic preparedness?



Sponsors

Frontier Science Foundation Hellas, Greece



Izmir Metropolitan Municipality, Türkiye



Chapman & Hall (CRC Press), UK



Visbanking, USA



StataCorp LLC, USA



TÜBİTAK, Türkiye





Awards

EMR 2023 Best Paper and Poster Awards

Awards will be given at the EMR 2023 for two best oral presentations and two best poster presentations by undergraduate, graduate students, and postdoctoral researchers. Entries will be judged on the quality of both the presentation and underlying research. The awards consist of a certificate and a prize.

To be eligible for this award:

- A student must not have defended her/his thesis nor completed her/his final degree requirements by the conference date.
- The postdoctoral researcher must not be working full-time at any company/institution.
- Indicate your best paper/poster application at the second page of the abstract template.

Travel Award (Student) by FSH-F

[Frontier Science Foundation Hellas](#) (FSF-H) is supporting student awards in honor of Prof. Lagakos as for the previous meetings. There will be 3 awards for students.

The awards include, up to 1000 Euros,

- air travel (economy class),
- registration fees, and
- accommodation for the days of the conference.

Those interested must apply by

- sending an extended abstract for up to 4 pages using the [Latex template](#),
- sending a letter from their supervisor that certifies that they are students (use as Subject "EMR Student Awards") to the address karlis@aueb.gr

Note that depending on the submissions the committee may conclude that one of the awards may be decided after the presentations in Izmir.



General Information

Conference Venue

The Conference will take place at the **Ege Palas Hotel**.

Alsancak, Cumhuriyet Blv No:210
35220 Konak/İzmir Türkiye
Tel: +90 (232) 4639090

Climate and Clothing

In May the weather is nice and warm during the whole day. Mean temperature ranges from 15-25 °C.

Country Dialing

Code +90

Electrical voltage

The voltage in Türkiye is 220 volt.

Liability and Insurance

The organizers have no responsibility whatsoever for injury or damage involving persons and property during the Conference. Participants are advised to carry their own personal insurance during their stay in Türkiye.

Name badges

All participants and accompanying persons must wear the Conference identification badge in a visible place. Entrance to meeting hall, poster and exhibition areas will not be permitted to any person without badge.

Official language

The official Language of the Conference will be English. Simultaneous translation will not be provided.

Tipping

For taxi and restaurants, the service charge is included in the price. You may add a tip at your own discretion to indicate appreciation of exceptionally good service.



Social Programme

After meeting at the hotel, we will depart for Ephesus Social Program. The first visit will be to the Virgin Mary House where she spent her last days, preferred to live in this remote place rather than in a crowded place.



Then, we will visit the most worth-seeing remains of the city such as The Marble Street, The Odeon, Bouleterion, The Temple of Hadrian, The Trajan's Fountain, The Mosaic street, The Agora, The Baths, The house, The Great Theatre, Harbour Street, and the third largest library of the Ancient World: The famous Celsus Library.



We follow our tour with the famous Sirince, a beautiful hill town. It is a fascinating village where you can feel traditional Turkish life. You can walk around, taste local wine, and stop for a cup of Turkish coffee. This tour is easy and relaxing and lets you take lovely pictures of different scenery and traditional Culture. After the Sirince tour, participants who prefers the Social Program (A) will be transferred to the city center. For the Social Program (B), we will have dinner at Arkas Vineyards (Lucien Arkas Bağları, <https://lawines.com.tr/>) accompanied by local wines. After dinner, our guests will be transferred to the city center.





Conference Program

May 8, 2023 (MONDAY)

TIME	BALO SALONU (1st FLOOR FLOOR M)	NAMIK SEVIK (4th FLOOR FLOOR P)	ACELYA (4th FLOOR FLOOR P)
09:00-10:30	COURSE I <i>Introduction to explainable machine learning with examples in healthcare</i> Przemyslaw Biecek	COURSE II <i>Removing unwanted variation from large-scale RNA sequencing data with PRPS</i> Ramyar Molania & Marie Trussart	COURSE III <i>How to use cloud technologies for reliable and responsible Data Science projects?</i> Erdal Coşgun, Vincent Carey & Deniz İlhan Topçu
10:30-11:00		Coffee break	
11:00-12:30	COURSE I <i>Introduction to explainable machine learning with examples in healthcare</i> Przemyslaw Biecek	COURSE II <i>Removing unwanted variation from large-scale RNA sequencing data with PRPS</i> Ramyar Molania & Marie Trussart	COURSE III <i>How to use cloud technologies for reliable and responsible Data Science projects?</i> Erdal Coşgun, Vincent Carey & Deniz İlhan Topçu
12:30-13:30		Lunch	
13:30-15:00	HONORARY MINI SYMPOSIUM: SESSION I <u>Chair: Yoav Benjamini</u> <i>Sparse kernel factor analysis model for high-dimensional undersampled cancer data sets and supervised classification with information complexity criterion</i> Hamparsum Bozdoğan <i>Variational multiple imputation in high-dimensional regression models with missing response</i> Recai M. Yücel GENESELECTML: A comprehensive way of gene selection for RNA-seq data via machine-learning algorithms Özlem İlk	COURSE II <i>Removing unwanted variation from large-scale RNA sequencing data with PRPS</i> Ramyar Molania & Marie Trussart	-
15:00-15:30		Coffee break	
15:30-17:30	HONORARY MINI SYMPOSIUM: SESSION II <u>Chair: Atilla Halil Elhan</u> <i>Confidence intervals for the Weitzman overlapping coefficient: the Binormal approach and alternatives</i> Benjamin Reiser <i>Mapping the disease prevalence of TÜRKİYE</i> Mehmet Kocak <i>The applied statistical (data) scientist in a high-profile and societal environment - IBS and EMR -</i> Geert Molenberghs Ergun Karaağaoğlu & Turkish Journal of Biochemistry Yahya Laleli	COURSE II <i>Removing unwanted variation from large-scale RNA sequencing data with PRPS</i> Ramyar Molania & Marie Trussart	-
17:30-18:30	HONORARY MINI SYMPOSIUM: SESSION III <u>Co-chairs: İlker Ünal & Erdal Coşgun</u> H. Refik Burgut A. Ergun Karaağaoğlu	-	-
18:30-20:30	Opening Cocktail (at Saint Voukolos Church) (How to Reach?: https://tinyurl.com/2bxrhfu3)		



May 9, 2023 (TUESDAY)

TIME	BALO SALONU (1st FLOOR FLOOR M)	NAMIK SEVIK (4th FLOOR FLOOR P)
08:00-08:30	Registration Opens	
	WELCOME SESSION	
08:30-09:30	<p>Talk by EMR 2023 Conference Chair (Assoc.Prof.Dr.Gökmen Zararsız)</p> <p>Talk by EMR-IBS President (Prof.Dr.Konstantinos Fokianos)</p> <p>Talk by TÜBİTAK President (Prof.Dr.Hasan Mandal)</p> <p>Recital of Turkish Music</p>	-
	KEYNOTE LECTURE	
09:30-10:15	<p>Chair: Karen Kafadar</p> <p>BIOCONDUCTOR: Evolving an open source ecosystem for genomic data science</p> <p>Vincent Carey</p>	-
10:15-10:30	Coffee Break	
	CONTRIBUTED SESSION 1: MULTIVARIATE STATISTICS	
	Co-chairs: Daniel Yekutieli & Özlem İlk Dağ	
	<i>Estimation of average causal effect in clustered data with covariate measurement error</i>	
	Recai M. Yücel, Raina E. Josberger & Meng Wu	
	<i>The chernoff faces method for visualizing complex data: An application for identifying differences between Covid-19 and control groups</i>	
	Elif Kavmaz, Ferhan Elmalı, Büşra Emir, Fatma Ezgi Can & Mustafa Agah Tekindal	
	<i>Viral load dynamics of Sars-Cov-2 delta and omicron variants following multiple vaccine doses and previous infection</i>	
	Naama M. Kopelman, Yonatan Woodbridge, Sharon Amit & Amit Huppert	
	<i>Evaluating univariate, multivariate reference interval methods: A comparative analysis</i>	
	Esra Kutsal Mergen & Sevilay Karahan	
	<i>Evaluation of objective structured examination tool with classical testing institution, generalizability theory and item response theory</i>	
	M. Yasemin Akşehırlı Sevfeli, Atilla Halil Elhan, Zeynep Baykan, Gözde E. Zararsız, Orhun Öztürk, Gökmen Zararsız & Ahmet Öztürk	
10:30-12:00	<p style="text-align: center;">INVITED SESSION 1: FUNCTIONAL DATA ANALYSIS</p> <p style="text-align: center;">Chair: Philip T. Reiss</p> <p><i>Functional data analyses to account for and interpret glucose patterns in continuous glucose monitoring</i></p> <p style="text-align: center;">Emrah Gecili & Rhonda Szczesniak</p> <p><i>Incorporating shared peptides for improved inference on the proteins' abundance based on mass spectrometry data</i></p> <p style="text-align: center;">Malgorzata Bogdan</p> <p><i>Functional additive models for shapes and forms of plane curves</i></p> <p style="text-align: center;">Almond Stöcker, Lisa Steyer & Sonja Greven</p>	-
12:00-12:15	Coffee Break	
	PLENARY LECTURE	
12:15-13:00	<p style="text-align: center;">Chair: H. Refik Burgut</p> <p><i>Replicability issues in medical research: Science and politics</i></p> <p style="text-align: center;">Yoav Benjamini</p>	-
13:00-14:00	Lunch	



May 9, 2023 (TUESDAY)

14:00-14:45	<p>MARVIN ZELEN MEMORIAL LECTURE</p> <p>Chair: <u>Urania Dafni</u></p> <p><i>Dynamic evaluation of Covid-19 clinical states by means of multi state models</i></p> <p>Guadalupe Gómez Melis</p>	-
14:45-16:15	<p>INVITED SESSION 2: BIOINFORMATICS AND HIGH DIMENSIONAL DATA ANALYSIS</p> <p>Chair: <u>Vincent Carey</u></p> <p><i>Removing unwanted variation from large gene expression data with RUV-III-PRPS</i></p> <p>Ramyar Molania</p> <p><i>Validation of model selection procedures in high-dimensional analysis</i></p> <p><u>Victor Kipnis</u>, Grant Izmirlian, Douglas Midthune</p> <p><i>Rapidly advancing CRISPR systems hold a great potential in research on drug targets</i></p> <p>Göknur Giner</p>	<p>YOUNG STATISTICIANS SHOWCASE</p> <p>Chair: <u>Dimitris Karlis</u></p> <p><i>Modified SHAP method for seasonal vaccination status</i></p> <p>Bekir Çetintav, Selim Çetin & <u>Ahmet Yalçın</u></p> <p><i>SimElegans: A Shiny app for GO analysis</i></p> <p><u>İrem Kahveci</u>, H. Furkan Kepenek & Dinçer Göksülük</p> <p><i>Real-time detection of the start and subsequent epidemic states of HIV outbreaks</i></p> <p><u>Valia Baralou</u>, Argiro Karakosta, Christos Thomadakis, Nikos Demiris, Nikos Pantazis, Olga Anagnostou, Christos Danopoulos, Dimitris Katsiris & Giota Touloumi</p> <p><i>Model based clustering for spatial data</i></p> <p>Anna Nalpantidi</p> <p><i>Multilevel Bayesian network to model child morbidity using Gibbs sampling</i></p> <p><u>Bezalem Eshetu Yirdaw</u> & Legesse Kassa Debusho</p>
16:15-16:45	Coffee Break	
16:45-17:30	<p>PLENARY LECTURE</p> <p>Chair: <u>Hamparsum Bozdoğan</u></p> <p><i>Simulation approach in the design and planning cancer screening trials</i></p> <p>Ping Hu</p>	-
17:30-18:30	<p>ROUND TABLE</p> <p>Co-chairs: <u>Geert Molenberghs & David M. Steinberg</u></p> <p><i>The role of statistical experts during the COVID-19 pandemic</i></p> <p>Arne Bathke</p> <p>Ralph Brinks</p> <p>Filomena Maggino</p> <p>Bhramar Mukherjee</p> <p>Amit Huppert</p>	<p>CONTRIBUTED SESSION 2: STATISTICAL MODELING AND COMPUTATIONAL METHODS</p> <p>Co-chairs: <u>George Dennis Papandonatos & Ferhan Elmalı</u></p> <p><i>Replicability across multiple studies</i></p> <p><u>Ruth Heller</u> & Marina Bogomolov</p> <p><i>Modelling longitudinal cognitive test data with ceiling effects and left skewness</i></p> <p><u>Denitsa Grigorova</u>, Dean Palejev & Ralitza Gueorguieva</p> <p><i>Branching modelling of mutations and risk assessment in cancer research</i></p> <p><u>Maroussia Slavtchova-Bojkova</u> & Kaloyan Vitanov</p> <p><i>Joint spatiotemporal modelling of human immunodeficiency virus and tuberculosis in ethiopia using a Bayesian hierarchical approach</i></p> <p><u>Legesse Kassa Debusho</u> & Leta Lencha Gemechu</p> <p><i>Comparing frequentist and Bayesian approaches for mixed design anova in repeated measurements: a simulation study with exponential distributions</i></p> <p><u>Zeynep Özel</u>, Ebru Kaya Başar & Mustafa Agah Tekindal</p>
18:30-19:00	POSTER SESSION	
19:00-19:45	EMR BUSINESS MEETING	
20:00-23:00	GALA Dinner (at North Pier's İzmir) (How to Reach?: https://tinyurl.com/yckkh8s3)	



May 10, 2023 (WEDNESDAY)

TIME	BALO SALONU (1st FLOOR FLOOR M)	NAMIK SEVIK (4th FLOOR FLOOR P)
08:30-10:00	<p>INVITED SESSION 3: CLINICAL BIOSTATISTICS AND SURVIVAL ANALYSIS Chair: David M. Steinberg <i>Missing time-dependent covariate values in a Cox model – joint models approach versus combination of multiple imputation and joint models</i> Havi Murad, Nirit Agay & Rachel Dankner <i>Estimating optimal individualized treatment rules with multistate processes</i> Giorgos Bakoyannis <i>Statistical inference for complex time-to-event data under non-randomized cohorts</i> KyungMann Kim, Tuo Wang, Lu Mao & Aldo Cocco</p>	<p>CONTRIBUTED SESSION 3: SURVIVAL ANALYSIS Co-chairs: Guadalupe GómezMelis & Victor Kipnis <i>Data-driven simulations for quantitative bias analyses in real-world survival analyses</i> Michał Abrahamowicz, Marie-Eve Beauchamp, Anne-Laure Boulesteix, Tim Morris, Willi Sauerbrei & Jay Kaufman <i>Conditional randomization test for average treatment effect with survival forest</i> Mehmet Ali Kaygusuz & Vilda Purutçuoğlu <i>Reconstructing survival data from published Kaplan-Meier curves</i> Georgia Rompoti, Dimitris Karlis & Urania Dafni <i>Methodological issues with proportional hazard models</i> Maria-Tereza Dellaporta, Dimitris Karlis & Urania Dafni</p>
10:00-10:30	Coffee Break	
10:30-11:15	<p>PLENARY LECTURE Chair: Yoav Benjamini <i>Accounting for overdiagnosis in estimating components of survival time in randomized cancer screening trials</i> Karen Kafadar</p>	-
11:15-12:45	<p>INVITED SESSION 4: RECENT ADVANCES IN BIOSTATISTICS Chair: Ruth Heller <i>Biostatistics and SARS-CoV-2: research, policy advice, and communication</i> Geert Molenberghs <i>In Silico: Simulators and emulators in the European human brain project</i> David M. Steinberg, Mira Marcus-Kalish, David Refaeli, Gilad Shapira & Ella Shaposhnik <i>Incorporating fuzzy logic philosophy in the evaluation of forecasting models</i> Çağdaş Hakan Aladağ</p>	<p>CONTRIBUTED SESSION 4: GENOMICS AND PROTEOMICS Co-chairs: Ramvar Molania & Przemyslaw Biecek <i>Stacking based approaches for survival analysis of RNA-sequencing data</i> Ahu Cephe, Necla Koçhan, Ahmet Sezgin, Gözde Ertürk Zararsız, Erdem Karabulut & Gökmen Zararsız <i>Robust protein co-expression network for Covid-19</i> Ayça Ölmez & Aylin Alın <i>Infoget4gene: A user-friendly web app for genetic data analysis using R shiny</i> Hamdi Furkan Kepenek, İrem Kahveci & Dinçer Göksülük <i>Parametric bootstrap based simulation on identification of differentially expressed genes: Which one of boruta or elastic net performs better?</i> Merve Kaşıkçı, Özgür Saman & Osman Dağ <i>Determination of intron retention in gastric cancer RNA-sequence data by IRFinder-S bioinformatics algorithm</i> Esmâ Gamze Aksel, Vahap Eldem, Selim Can Kuralay & Gökmen Zararsız</p>
13:00-18:00	SOCIAL Program (Ephesus & Mary's House & Sirince Village)	
18:00-21:00	SOCIAL Program Dinner (at Lucien Arkas Vineyards)	



May 11, 2023 (THURSDAY)

TIME	BALO SALONU (1st FLOOR FLOOR M)	NAMIK SEVIK (4th FLOOR FLOOR P)
08:30-10:00	<p>CONTRIBUTED SESSION 5A: MACHINE LEARNING AND ARTIFICIAL INTELLIGENCE - I Co-chairs: Aris Perperoglou & Mustafa Cavus <i>Feature extraction and biomarker analysis for differentiating colon polyps from colonoscopic images</i> Refika Sultan Doğan, Ebru Aker, Serkan Doğan & Bülent Yılmaz <i>WRSmoonRF: Weighted robust sufficient M-out-of-N regression forest</i> Aylin Alın <i>Optimizing number of hidden layer and hyper-parameters of deep neural network by Bayesian optimization</i> Yasin Gürmez, Duygu Korkmaz Yalçın & Sıddık Keskin <i>Automated machine learning approach in clinical settings: Predicting the future risks</i> Didem Turgut, Deniz İlhan Topçu, Samet Şenel & Cüneyt Özden <i>A comprehensive comparison of low density lipoprotein cholesterol equations</i> Serra İlayda Yerlitas, Gözde Ertürk Zararsız, Halef Okan Doğan, Serkan Bolat, Necla Koçhan, Ahu Cephe, Gökmen Zararsız & Arrigo F.G. Cicero</p>	<p>CONTRIBUTED SESSION 5B: CLUSTERING AND CLASSIFICATION Co-chairs: Malgorzata Bogdan & Osman Dağ <i>Bioinformatics and biostatistical models for analysis and prognosis of antimicrobial resistance</i> Maya Zhelyazkova, Stefan Tsonev & Dimitar Vassilev <i>Deep neural networks for average treatment effect on biological networks</i> Mehmet Ali Kaygusuz & Vilda Purutçuoğlu <i>Simultaneous scoring of clusters in recursive cluster elimination, applied on transcriptomic data analysis</i> Nurten Bulut, Burcu Bakır-Güngör, Bahjat F. Qaqish & Malik Yousef <i>The G-S-M, grouping, scoring and modeling approach. Application of biological domain knowledge for groups selection on gene expression data</i> Malik Yousef & Burcu Bakır-Güngör <i>The effect of missing data imputation methods on classification performance according to different missing rates in high dimensional data</i> Buğra Varol, İmran Kurt Ömürlü & Mevlüt Türe</p>
10:00-10:30	Coffee Break	
10:30-11:15	<p>PLENARY LECTURE Chair: KyungMann Kim <i>20 years of statistical bioinformatics: a brief history of the limma and edgeR packages</i> Gordon Smyth</p>	-
11:15-12:45	<p>INVITED SESSION 5: BRAIN IMAGING METHODS Chair: R. Todd Ogden <i>Nonparametric functional data modeling of pharmacokinetic processes with applications in dynamic pet imaging</i> R. Todd Ogden <i>Investigating directional causality in multichannel brain signals: Threshold autoregressive modeling based approach</i> Sipan Aslan & Hernando Ombao <i>How reliable is resting-state functional connectivity?</i> Xu Meng, Philip T. Reiss & Ivor Cribben</p>	<p>CONTRIBUTED SESSION 6: MACHINE LEARNING AND ARTIFICIAL INTELLIGENCE - II Co-chairs: Gökmen Zararsız & Necla Koçhan <i>Health space model using deep learning</i> Taesung Park & Chanhee Lee <i>How to explain carbohydrate metabolism disorders using machine learning models?</i> Deniz İlhan Topçu & Banu İşbilen Başok <i>Abnormality detection and classification in mammography images via convolutional neural networks</i> Hanife Avcı, Gamze Durhan, Figen Demirkazık, Meltem Gülsün Akpınar & Jale Karakaya <i>Estimating the cost of Covid-19 to Turkish tourism with time series and machine learning models</i> Günal Bilek <i>Assessing prediction accuracy of joint models: A novel approach based on mutual information criterion</i> Merve Başol Gökşülük, Dincer Gökşülük & A. Ergun Karaağaoğlu</p>
12:45-13:45	Lunch	
13:45-14:30	<p>PLENARY LECTURE Chair: A. Ergun Karaağaoğlu <i>Towards reliable empirical evidence in methodological biostatistical research: recent developments and remaining challenges</i> Anne-Laure Boulesteix</p>	-
14:30-16:15	<p>INVITED SESSION 6: MACHINE-LEARNING AND COMPUTER SCIENCE Chair: Taesung Park <i>Enhancing clinical trial power through covariate adjustment models: An investigation of methods, claims and realistic expectations</i> Aris Perperoglou <i>Cardiac classification using machine learning models with information complexity</i> Hamparsum Bozdogan <i>Bringing group sparsity to Bayesian hierarchical models: multivariate Bernoulli distribution is here to help</i> Zeynep Bal, Gülay Başarır, Mehmet Gönen</p>	<p>CONTRIBUTED SESSION 7: CLINICAL TRIALS AND DIAGNOSTIC TESTS Co-chairs: Aylin Alın & Sevilay Karahan <i>Practical considerations for statistical models and their implementations of phase I dose-escalation oncology trials</i> Burak Kürsad Günhan, Pavel Mozgunov & Anja Victor <i>A review on randomized controlled trials in emergency medicine: Methodological issues</i> Demet Arı, Pınar Günel, Buket İpek Berk, İhsan Berk & Vildan Sümbüloğlu <i>Nonparametric estimation of distribution function using ranked set sampling with unequal probabilities</i> Yusuf Can Sevil & Tuğba Özal Yıldız <i>Public health-focused use of Covid-19 rapid antigen PCR tests</i> Yonatan Woodbridge, Yair Goldberg, Sharon Amit, Naama M. Kopelman, Micha Mandel & Amit Huppert <i>Diabetic retinopathy diagnosis and classification</i> Berk Piskin, Aylin Alın, Rim Khazin, Ahmet Mert Saygu & Ahmet Ömer Özgür <i>A new estimator for the discrimination accuracy in a four-class classification problem</i> Elena Nardi</p>
16:15-16:45	Coffee Break	
16:45-17:15	Closing Remarks & Award Ceremony	



TABLE OF CONTENT

KEYNOTE LECTURE	1
BIOCONDUCTOR: EVOLVING AN OPEN SOURCE ECOSYSTEM FOR GENOMIC DATA SCIENCE	2
<i>Vincent Carey</i> ¹	2
MARVIN ZELEN MEMORIAL LECTURE	3
DYNAMIC EVALUATION OF COVID-19 CLINICAL STATES BY MEANS OF MULTI STATE MODELS.....	4
<i>Guadalupe Gómez Melis</i> ¹	4
PLENARY LECTURES	5
TOWARDS RELIABLE EMPIRICAL EVIDENCE IN METHODOLOGICAL BIOSTATISTICAL RESEARCH: RECENT DEVELOPMENTS AND REMAINING CHALLENGES.....	6
<i>Anne-Laure Boulesteix</i> ¹	6
20 YEARS OF STATISTICAL BIOINFORMATICS: A BRIEF HISTORY OF THE LIMMA AND edgeR PACKAGES	7
<i>Gordon Smyth</i> ¹	7
ACCOUNTING FOR OVERDIAGNOSIS IN ESTIMATING COMPONENTS OF SURVIVAL TIME IN RANDOMIZED CANCER SCREENING TRIALS	8
<i>Karen Kafadar</i> ¹	8
SIMULATION APPROACH IN THE DESIGN AND PLANNING CANCER SCREENING TRIALS	9
<i>Ping Hu</i> ¹	9
REPLICABILITY ISSUES IN MEDICAL RESEARCH: SCIENCE AND POLITICS	10
<i>Yoav Benjamini</i> ¹	10
INVITED TALKS	11
FUNCTIONAL ADDITIVE MODELS FOR SHAPES AND FORMS OF PLANE CURVES	12
<i>Almond Stöcker</i> ¹ , <i>Lisa Steyer</i> ² , <i>Sonja Greven</i> ³	12
ENHANCING CLINICAL TRIAL POWER THROUGH COVARIATE ADJUSTMENT MODELS: AN INVESTIGATION OF METHODS, CLAIMS AND REALISTIC EXPECTATIONS.....	13
<i>Aris Perperoglou</i> ¹	13
INCORPORATING FUZZY LOGIC PHILOSOPHY IN THE EVALUATION OF FORECASTING MODELS	14
<i>Cagdas Hakan Aladag</i> ¹	14
IN SILICO: SIMULATORS AND EMULATORS IN THE EUROPEAN HUMAN BRAIN PROJECT	15
<i>David M. Steinberg</i> ¹ , <i>Mira Marcus-Kalish</i> ¹ , <i>David Refaeli</i> ¹ , <i>Gilad Shapira</i> ¹ , <i>Ella Shaposhnik</i> ¹	15
FUNCTIONAL DATA ANALYSES TO ACCOUNT FOR AND INTERPRET GLUCOSE PATTERNS IN CONTINUOUS GLUCOSE MONITORING	16
<i>Emrah Gecili</i> ^{1, 2} , <i>Rhonda Szczesniak</i> ^{1, 2}	16
BIOSTATISTICS AND SARS-COV-2: RESEARCH, POLICY ADVICE, AND COMMUNICATION	17
<i>Geert Molenberghs</i> ¹	17
ESTIMATING OPTIMAL INDIVIDUALIZED TREATMENT RULES WITH MULTISTATE PROCESSES.....	18
<i>Giorgos Bakoyannis</i> ¹	18
RAPIDLY ADVANCING CRISPR SYSTEMS HOLD A GREAT POTENTIAL IN RESEARCH ON DRUG TARGETS	19
<i>Göknur Giner</i> ¹	19
SPARSE KERNEL FACTOR ANALYSIS MODEL FOR HIGH-DIMENSIONAL UNDERSAMPLED CANCER DATA SETS AND SUPERVISED CLASSIFICATION WITH INFORMATION COMPLEXITY CRITERION	20
<i>Hamparsum Bozdogan</i> ¹	20
MISSING TIME-DEPENDENT COVARIATE VALUES IN A COX MODEL – JOINT MODELS APPROACH VERSUS COMBINATION OF MULTIPLE IMPUTATION AND JOINT MODELS	21
<i>Havi Murad</i> ¹ , <i>Nirit Agay</i> ^{1, 2} , <i>Rachel Dankner</i> ^{2, 3}	21
STATISTICAL INFERENCE FOR COMPLEX TIME-TO-EVENT DATA UNDER NON-RANDOMIZED COHORTS.....	23
<i>KyungMann Kim</i> ¹ , <i>Tuo Wang</i> ¹ , <i>Lu Mao</i> ¹ , <i>Aldo Cocco</i> ²	23
INCORPORATING SHARED PEPTIDES FOR IMPROVED INFERENCE ON THE PROTEINS' ABUNDANCE BASED ON MASS SPECTROMETRY DATA	24
<i>Malgorzata Bogdan</i> ¹	24
BRINGING GROUP SPARSITY TO BAYESIAN HIERARCHICAL MODELS: MULTIVARIATE BERNOULLI DISTRIBUTION IS HERE TO HELP	25
<i>Mehmet Gönen</i> ¹	25



HOW RELIABLE IS RESTING-STATE FUNCTIONAL CONNECTIVITY?	26
<i>Xu Meng¹, Philip T. Reiss¹, Ivor Cribben²</i>	26
NONPARAMETRIC FUNCTIONAL DATA MODELING OF PHARMACOKINETIC PROCESSES WITH APPLICATIONS IN DYNAMIC PET IMAGING	27
<i>R. Todd Ogden¹</i>	27
REMOVING UNWANTED VARIATION FROM LARGE GENE EXPRESSION DATA WITH RUV-III-PRPS.....	28
<i>Ramyar Molania¹</i>	28
INVESTIGATING DIRECTIONAL CAUSALITY IN MULTICHANNEL BRAIN SIGNALS: THRESHOLD AUTOREGRESSIVE MODELING BASED APPROACH.....	29
<i>Sipan Aslan^{1,3}, Hernando Ombao²</i>	29
VALIDATION OF MODEL SELECTION PROCEDURES IN HIGH-DIMENSIONAL ANALYSIS.....	30
<i>Victor Kipnis¹, Grant Izmirlian¹, Douglas Midthune¹</i>	30
SYMPOSIUM TALKS	31
CONFIDENCE INTERVALS FOR THE WEITZMAN OVERLAPPING COEFFICIENT: THE BINORMAL APPROACH AND ALTERNATIVES.....	32
<i>Benjamin Reiser¹</i>	32
THE APPLIED STATISTICAL (DATA) SCIENTIST IN A HIGH-PROFILE AND SOCIETAL ENVIRONMENT	33
- IBS AND EMR -	33
<i>Geert Molenberghs¹</i>	33
CARDIAC CLASSIFICATION USING MACHINE LEARNING MODELS WITH INFORMATION COMPLEXITY.....	34
<i>Hamparsum Bozdogan¹</i>	34
MAPPING THE DISEASE PREVALENCE OF TURKIYE.....	35
<i>Mehmet Koçak¹</i>	35
GENESECTML: A COMPREHENSIVE WAY OF GENE SELECTION FOR RNA-SEQ DATA VIA MACHINE LEARNING ALGORITHMS	36
<i>Ozlem Ilk¹</i>	36
VARIATIONAL MULTIPLE IMPUTATION IN HIGH-DIMENSIONAL REGRESSION MODELS WITH MISSING RESPONSES.....	37
<i>Recai M. Yücel¹</i>	37
ERGUN KARAAĞAOĞLU & TURKISH JOURNAL OF BIOCHEMISTRY	38
<i>Yahya Laleli¹</i>	38
ORAL PRESENTATIONS.....	39
OP1. ESTIMATION OF AVERAGE CAUSAL EFFECT IN CLUSTERED DATA WITH COVARIATE MEASUREMENT ERROR	40
<i>Recai M. Yucel¹, Raina E. Josberger², Meng Wu³</i>	40
OP2. VIRAL LOAD DYNAMICS OF SARS-COV-2 DELTA AND OMICRON VARIANTS FOLLOWING MULTIPLE VACCINE DOSES AND PREVIOUS INFECTION	41
<i>Naama M. Kopelman¹, Yonatan Woodbridge^{1,2}, Sharon Amit³, Amit Huppert^{2,4}</i>	41
OP3. EVALUATING UNIVARIATE, MULTIVARIATE REFERENCE INTERVAL METHODS: A COMPARATIVE ANALYSIS	42
<i>Esra Kutsal Mergen¹, Sevilay Karahan²</i>	42
OP4. EVALUATION OF OBJECTIVE STRUCTURED EXAMINATION TOOL WITH CLASSICAL TESTING INSTITUTION, GENERALIZABILITY THEORY AND ITEM RESPONSE THEORY.....	43
<i>M.Yasemin Akşehirli Seyfeli¹, Atilla H. Elhan², Zeynep Baykan³,</i>	43
<i>Gözde Ertürk Zararsız⁴, Orhun Öztürk⁵, Gökmen Zararsız⁴, Ahmet Öztürk⁴</i>	43
OP5. REPLICABILITY ACROSS MULTIPLE STUDIES.....	44
<i>Ruth Heller¹, Marina Bogomolov²</i>	44
OP6. MODELLING LONGITUDINAL COGNITIVE TEST DATA WITH CEILING EFFECTS AND LEFT SKEWNESS	45
<i>Denitsa Grigorova¹, Dean Palejev², Ralitzia Gueorguieva³</i>	45
OP7. BRANCHING MODELLING OF MUTATIONS AND RISK ASSESSMENT IN CANCER RESEARCH.....	46
<i>Maroussia Slavtchova-Bojkova¹, Kaloyan Vitanov²</i>	46
OP8. JOINT SPATIOTEMPORAL MODELLING OF HUMAN IMMUNODEFICIENCY VIRUS AND TUBERCULOSIS IN ETHIOPIA USING A BAYESIAN HIERARCHICAL APPROACH	47
<i>Legesse Kassa Debusho¹, Leta Lencha Gemechu²</i>	47
OP9. COMPARING FREQUENTIST AND BAYESIAN APPROACHES FOR MIXED DESIGN ANOVA IN REPEATED MEASUREMENTS: A SIMULATION STUDY WITH EXPONENTIAL DISTRIBUTIONS	48
<i>Zeynep Özel¹, Ebru Kaya Başar², Mustafa Agah Tekindal¹</i>	48
OP10. DATA-DRIVEN SIMULATIONS FOR QUANTITATIVE BIAS ANALYSES IN REAL-WORLD SURVIVAL ANALYSES	49



<i>Michal Abrahamowicz¹, Marie-Eve Beauchamp¹, Anne-Laure Boulesteix²,</i>	<i>49</i>
<i>Tim Morris³, Willi Sauerbrei⁴, Jay Kaufman¹</i>	<i>49</i>
OP11. CONDITIONAL RANDOMIZATION TEST FOR AVERAGE TREATMENT EFFECT WITH SURVIVAL FOREST	50
<i>Mehmet Ali Kaygusuz¹, Vilda Purutçuoğlu</i>	<i>50</i>
OP12. RECONSTRUCTING SURVIVAL DATA FROM PUBLISHED KAPLAN-MEIER CURVES	51
<i>Georgia Rompoti^{1,2}, Dimitris Karlis³, Urania Dafni^{1,2}.....</i>	<i>51</i>
OP13. METHODOLOGICAL ISSUES WITH PROPORTIONAL HAZARD MODELS	52
<i>Maria-Tereza Dellaporta^{1,2}, Dimitris Karlis¹, Urania Dafni^{2,3}.....</i>	<i>52</i>
OP14. STACKING BASED APPROACHES FOR SURVIVAL ANALYSIS OF RNA-SEQUENCING DATA	53
<i>Ahu Cephe^{1,2}, Necla Koçhan³, Ahmet Sezgin⁴, Gözde Ertürk Zararsız^{2,5}, Erdem Karabulut⁶, Gökmen Zararsız^{2,5}</i>	<i>53</i>
OP15. ROBUST PROTEIN CO-EXPRESSION NETWORK FOR COVID-19	54
<i>Ayca Olmez¹, Aylin Alin²</i>	<i>54</i>
OP16. infoget4gene: A USER-FRIENDLY WEB APP FOR GENETIC DATA ANALYSIS USING R SHINY	55
<i>Hamdi Furkan Kepenek¹, İrem Kahveci¹, Dinçer Goksülük²</i>	<i>55</i>
OP17. PARAMETRIC BOOTSTRAP BASED SIMULATION ON IDENTIFICATION OF DIFFERENTIALLY EXPRESSED GENES: WHICH ONE OF BORUTA OR ELASTIC NET PERFORMS BETTER?	56
<i>Merve Kasikci¹, Ozgur Saman¹, Osman Dag¹</i>	<i>56</i>
OP18. DETERMINATION OF INTRON RETENTION IN GASTRIC CANCER RNA-SEQUENCE DATA BY IRFINDER-S BIOINFORMATICS ALGORITHM.....	57
<i>Esmâ Gamze Aksel¹, Vahap Eldem², Selim Can Kuralay², Gökmen Zararsız³.....</i>	<i>57</i>
OP19. FEATURE EXTRACTION AND BIOMARKER ANALYSIS FOR DIFFERENTIATING COLON POLYPS FROM COLONOSCOPIC IMAGES.....	58
<i>Refika Sultan Doğan¹, Ebru Aker², Serkan Doğan³, Bülent Yılmaz⁴</i>	<i>58</i>
OP20. WRSmoonRF: WEIGHTED ROBUST SUFFICIENT M-OUT-OF-N REGRESSION FOREST	59
<i>Aylin Alin¹.....</i>	<i>59</i>
OP21. OPTIMIZING NUMBER OF HIDDEN LAYER AND HYPER-PARAMETERS OF DEEP NEURAL NETWORK BY BAYESIAN OPTIMIZATION...60	60
<i>Yasin Görmez¹, Duygu Korkmaz Yalçın², Siddik Keskin².....</i>	<i>60</i>
OP22. AUTOMATED MACHINE LEARNING APPROACH IN CLINICAL SETTINGS: PREDICTING THE FUTURE RISKS	61
<i>Didem Turgut^{1,2}, Deniz İlhan Topcu³, Samet Senel⁴, Cuneyt Ozden⁴</i>	<i>61</i>
OP23. BIOINFORMATICS AND BIOSTATISTICAL MODELS FOR ANALYSIS AND PROGNOSIS OF ANTIMICROBIAL RESISTANCE	62
<i>Maya Zhelyazkova¹, Stefan Tsonev², Dimitar Vassilev³.....</i>	<i>62</i>
OP24. DEEP NEURAL NETWORKS FOR AVERAGE TREATMENT EFFECT ON BIOLOGICAL NETWORKS.....	63
<i>Mehmet Ali Kaygusuz¹, Vilda Purutçuoğlu</i>	<i>63</i>
OP25. SIMULTANEOUS SCORING OF CLUSTERS IN RECURSIVE CLUSTER ELIMINATION, APPLIED ON TRANSCRIPTOMIC DATA ANALYSIS..64	64
<i>Nurten Bulut¹, Burcu Bakir-Gungor², Bahjat F. Qaqish³, Malik Yousef⁴.....</i>	<i>64</i>
OP26. THE G-S-M, GROUPING, SCORING AND MODELING APPROACH. APPLICATION OF BIOLOGICAL DOMAIN KNOWLEDGE FOR GROUPS SELECTION ON GENE EXPRESSION DATA.....	65
<i>Malik Yousef¹, Burcu Bakir-Gungor².....</i>	<i>65</i>
OP27. THE EFFECT OF MISSING DATA IMPUTATION METHODS ON CLASSIFICATION PERFORMANCE ACCORDING TO DIFFERENT MISSING RATES IN HIGH DIMENSIONAL DATA	66
<i>Buğra Varol¹, İmran Kurt Ömürlü², Mevlüt Türe².....</i>	<i>66</i>
OP28. HEALTH SPACE MODEL USING DEEP LEARNING	67
<i>Taesung Park¹, Chanhee Lee².....</i>	<i>67</i>
OP29. HOW TO EXPLAIN CARBOHYDRATE METABOLISM DISORDERS USING MACHINE LEARNING MODELS?	68
<i>Deniz İlhan Topcu¹, Banu Isbilen Basok¹.....</i>	<i>68</i>
OP30. ABNORMALITY DETECTION AND CLASSIFICATION ON MAMMOGRAPHY IMAGES VIA CONVOLUTIONAL NEURAL NETWORKS	69
<i>Hanife Avci¹, Gamze Durhan², Figen Demirkazık², Meltem Gülsün Akpınar², Jale Karakaya¹</i>	<i>69</i>
OP31. ESTIMATING THE COST OF COVID-19 TO TURKISH TOURISM WITH TIME SERIES AND MACHINE LEARNING MODELS	70
<i>Günal Bilek¹.....</i>	<i>70</i>
OP32. ASSESSING PREDICTION ACCURACY OF JOINT MODELS: A NOVEL APPROACH BASED ON MUTUAL INFORMATION CRITERION	71
<i>Merve Basol Goksuluk¹, Dincer Goksuluk¹, A. Ergun Karaagaoglu²</i>	<i>71</i>
OP33. PRACTICAL CONSIDERATIONS FOR STATISTICAL MODELS AND THEIR IMPLEMENTATIONS OF PHASE I DOSE-ESCALATION ONCOLOGY TRIALS	72
<i>Burak Kürsad Günhan¹, Pavel Mozgunov², Anja Victor¹.....</i>	<i>72</i>



OP34. THE CHERNOFF FACES METHOD FOR VISUALIZING COMPLEX DATA: AN APPLICATION FOR IDENTIFYING DIFFERENCES BETWEEN COVID-19 AND CONTROL GROUPS	73
<i>Elif Kaymaz¹, Ferhan Elmalı¹, Büşra Emir¹, Fatma Ezgi Can¹, Mustafa Ağah Tekindal¹</i>	73
OP35. A REVIEW ON RANDOMIZED CONTROLLED TRIALS IN EMERGENCY MEDICINE: METHODOLOGICAL ISSUES	74
<i>Demet Arı¹, Pınar Günel², Buket İpek Berk³, İhsan Berk², Vildan Sümbüloğlu²</i>	74
OP36. NONPARAMETRIC ESTIMATION OF DISTRIBUTION FUNCTION USING RANKED SET SAMPLING WITH UNEQUAL PROBABILITIES	75
<i>Yusuf Can Sevil¹, Tuğba Özkal Yıldız²</i>	75
OP37. A COMPREHENSIVE COMPARISON OF LOW DENSITY LIPOPROTEIN CHOLESTEROL EQUATIONS	76
<i>Serra İlayda Yerlitaş^{1,2}, Gözde Ertürk Zararsız^{1,2}, Halef Okan Doğan³, Serkan Bolat³,</i>	76
<i>Necla Koçhan⁴, Ahu Cephe⁵, Gökmen Zararsız^{1,2}, Arrigo F.G. Cicero^{6,7}</i>	76
OP38. PUBLIC HEALTH-FOCUSED USE OF COVID-19 RAPID ANTIGEN PCR TESTS	77
<i>Yonatan Woodbridge^{1,2}, Yair Goldberg³, Sharon Amit⁴, Naama M. Kopelman²,</i>	77
<i>Micha Mandel⁵ and Amit Huppert^{1,6}</i>	77
OP39. DIABETIC RETINOPATHY DIAGNOSIS AND CLASSIFICATION	78
<i>Berk Pişkin¹, Aylin Alın², Rim Khazhin³, Ahmet Mert Saygu⁴, Ahmet Ömer Özgür⁵</i>	78
OP40. A NEW ESTIMATOR FOR THE DISCRIMINATION ACCURACY IN A FOUR-CLASS CLASSIFICATION PROBLEM.....	79
<i>Elena Nardi¹</i>	79
OP41. REVIEW: ACCESS CONTROL IN ELECTRONIC MEDICAL RECORDS (EMR)	80
<i>Hilal Alnafisah¹, Rawaby Alsaaid¹, Fatimah M. Alturkistani¹</i>	80
POSTER PRESENTATIONS	81
PP1. FAST AND APPROXIMATE INFERENCE OF MULTILEVEL THRESHOLD AUTOREGRESSIVE MODEL FOR INTENSIVE LONGITUDINAL DATA VIA MEAN FIELD VARIATIONAL BAYES.....	82
<i>Azizur Rahman^{1,2}, Depeng Jiang²</i>	82
PP2. ARTIFICIAL INTELLIGENCE-BASED MORPHOLOGY ANALYSIS SYSTEM FOR BRAIN ORGANIDS	83
<i>Elifsu Polatlı^{1,2}, Burak Kahveci^{1,2}, Sinan Güven^{1,2,3}</i>	83
PP3. TRANSCRIPTOMIC PROFILING OF INDUCED PLURIPOTENT STEM CELL DERIVED LACRIMAL ORGANIDS	84
<i>Burak Kahveci^{1,2}, Gamze Koçak^{1,2}, Canan Aslı Utine^{1,3}, Adil Mardinoğlu^{4,5},</i>	84
<i>Gökhan Karakulah^{1,2}, Sinan Güven^{1,2,6}</i>	84
PP4. PREDICTA: QUICK AND ACCURATE TRIAGE TOOL	85
<i>Vahide Gül Türkmen¹, Burak Kahveci^{2,3}</i>	85
PP5. MODIFIED CLINICAL KERNEL USING A COX MODEL	86
<i>Seungyeoun Lee¹, Inyoung Kim², Hyunjae Lee¹</i>	86
PP6. NET BENEFIT IN CLINICAL DECISION MAKING PROCESS.....	87
<i>Duygu Korkmaz Yalçın^{1,2}, İlker Ünal¹</i>	87
PP7. ESTIMATION OF LOW-DENSITY LIPOPROTEIN CHOLESTEROL USING MACHINE LEARNING MODELS.....	88
<i>Necla Koçhan¹</i>	88
PP8. EVALUATION OF SURVIVAL TREE RANDOM SURVIVAL FOREST AND COX PROPORTIONAL HAZARD MODELS.....	89
<i>Duygu Korkmaz Yalçın¹, Siddik Keskin¹</i>	89
PP9. ALZHEIMER DISEASE CLASSIFICATION WITH MACHINE LEARNING METHOD	90
<i>Fatma Gül Gezer¹, Berfu Parçalı², Kevser Setenay Öner², Fezan Mutlu²</i>	90
PP10. A COMPREHENSIVE R SHINY WEB TOOL THAT COMBINES TWO CONTINUOUS DIAGNOSTIC TESTS.....	91
<i>Serra İlayda Yerlitaş^{1,2}, Serra Bersan Gengeç², Gözde Ertürk Zararsız^{1,2}</i>	91
<i>Selçuk Korkmaz³, Gökmen Zararsız^{1,2}</i>	91
PP11. A CIRCULAR HEATMAP VISUALIZATION APPROACH FOR INTERLABORATORY COMPARISONS IN RING STUDIES.....	92
<i>Gözde Ertürk Zararsız^{1,2}, Alexander Cecil³, Jutta Lintelmann², Gernot Poschet⁴, Jennifer Kirwan⁵,</i>	92
<i>Sven Schuchardt⁶, Xue Li Guan⁷, Daisuke Saigusa⁸, David Wishart⁹, Jiamin Zheng¹⁰, Rupasri Mandal¹⁰,</i>	92
<i>Lisa St. John-Williams¹¹, Kendra Adams¹¹, J. Will Thompson¹¹, Michael P. Snyder¹², Kevin Contrepois¹²,</i>	92
<i>Songlie Chen¹², Nadia Ashrafi¹³, Sumeyya Akyol¹³, Ali Yilmaz¹³, Stewart Graham¹³, Thomas M. O'Connell¹⁴,</i>	92
<i>Karl Kalecký^{15, 16}, Teodoro Bottiglieri^{15,1}, Tuan Hai Pham¹⁷; Therese Koal¹⁷, Jerzy Adamski^{18,19,20}, Gabi Kastenmüller²</i>	92
PP12. A META-ANALYSIS STUDY FOR DIAGNOSING SKIN CANCER WITH MACHINE LEARNING TECHNIQUES	94
<i>Gözde Ertürk Zararsız^{1,2}, Elif Çelik Gürbulak^{2,3}, Serra İlayda Yerlitaş^{1,2}, Selen Yılmaz Işıkan⁴, Abdullah Demirbaş⁵, İrem Eroğlu^{2,6}, Aleyna Erakçaoğlu^{2,6}, Ragıp Ertaş⁷, Ömer Faruk Elmas⁸, Gökmen Zararsız^{1,2}</i>	94
PP13. ESTABLISHING CONTINUOUS REFERENCE INTERVALS FOR THYROID FUNCTION TESTS	95



<i>Funda İpekten</i> ^{1,2,3} , <i>Gözde Ertürk Zararsız</i> ^{1,2} , <i>Halef Okan Doğan</i> ⁴ , <i>Çiğdem Karakükçü</i> ⁵ , <i>Gökmen Zararsız</i> ^{1,2}	95
PP14. ON PERFORMANCES OF DIFFERENT CORRELATION COEFFICIENTS	96
<i>Yeşim Uzun Uğur</i> ¹ , <i>Mehmet Mendeş</i> ¹	96
YOUNG STATISTICIANS SHOWCASE PRESENTATIONS	97
MODIFIED SHAP METHOD FOR SEASONAL VACCINATION STATUS	98
<i>Ahmet Yalcin</i> ¹ , <i>Bekir Cetintav</i> ² , <i>Selim Cetin</i> ³	98
SHINY APP FOR GO ANALYSIS (SimElegans)	99
<i>İrem Kahveci</i> ¹ , <i>Hamdi Furkan Kepenek</i> ¹ , <i>Diñçer Göksülük</i> ²	99
REAL-TIME DETECTION OF THE START AND SUBSEQUENT EPIDEMIC STATES OF HIV OUTBREAKS AMONG PEOPLE WHO INJECT DRUGS: INSIGHTS FROM FOUR EUROPEAN COUNTRIES	100
<i>Valia Baralou</i> ^{1*} , <i>Argiro Karakosta</i> ^{1*} , <i>Christos Thomadakis</i> ¹ , <i>Nikos Demiris</i> ² , <i>Nikos Pantazis</i> ¹ , <i>Olga Anagnostou</i> ³ , <i>Christos Danopoulos</i> ³ , <i>Dimitris Katsiris</i> ⁴ , <i>Giota Touloumi</i> ¹	100
MODEL BASED CLUSTERING FOR SPATIAL DATA	101
<i>Anna Nalpantidi</i> ¹	101
MULTILEVEL BAYESIAN NETWORK TO MODEL CHILD MORBIDITY USING GIBBS SAMPLING	102
<i>Bezalem Eshetu Yirdaw</i> ¹ , <i>Legesse Kassa Debusho</i> ¹	102
INDEX	103

KEYNOTE LECTURE

BIOCONDUCTOR: EVOLVING AN OPEN SOURCE ECOSYSTEM FOR GENOMIC DATA SCIENCE

Vincent Carey¹

¹Harvard Medical School, USA

e-mail: *stvjc@channing.harvard.edu*

The Bioconductor project has coordinated the development and distribution of software, data, documentation and training resources for genomic data science for over 20 years. It seems likely that software from the project is on the desktop or laptop of the majority of bioinformaticians on the planet. In this talk I will review a number of core principles and methods of the project that have been important for creating a community of collaborative developers of resilient and performant analytical tools for an ever-expanding range of genome-scale assays. I will provide examples of evolving approaches of the project to addressing problems connected with using new reference genomes, large scale genome-wide association studies, and deep learning with single-cell multiomics experiments. Coupled to these examples will be illustrations of inter-language interfacing, profiling of resource usage, and development and deployment of inclusive teaching resources.

MARVIN ZELEN MEMORIAL LECTURE

DYNAMIC EVALUATION OF COVID-19 CLINICAL STATES BY MEANS OF MULTI STATE MODELS

Guadalupe Gómez Melis¹

¹*GRBIO: Research Group in Biostatistics and Bioinformatics
Departament d'Estadística i Investigació Operativa,
Universitat Politècnica de Catalunya-BarcelonaTech*

e-mail: lupe.gomez@upc.edu

Modeling the disease course of a hospitalized person regarding serious events is of great clinical relevance. Besides death, other clinically intermediate events such as admission to ICU or need of respiratory aid are important to identify prognostic factors and for clinical management. In this talk I will present the experience of the DIVINE multidisciplinary research team that was set up to achieve a deeper understanding of the severe form of the disease caused by the SARS-CoV-2 virus. The study was based on more than 4,000 subjects from several waves of the pandemic in five hospitals in the Barcelona metropolitan south region. Our project had the following four main goals: 1. Identification of clinically relevant prognostic factors for severe pneumonia, invasive respiratory support, death and discharge. The integrated combination of all these events was carried out through statistical multi-state models under Markov, semi-Markov and second order Markov assumptions; 2. A friendly interactive prediction web app, MSMPred, has been built to fit a generic multistate model. This app is amenable for diseases with different stages of severity; 3. Clustering analysis has been conducted to identify relevant hidden patient profiles which are associated with common sociodemographic and health profiles. This will help to design personalized treatments or to take a particular policy decision. Firstly, determining the number and the composition of the clusters and, secondly, comparing the patient's profiles among clusters; 4. 5th wave COVID-19 incubation time period has been estimated based on date, or a range of dates, where the infection might have occurred as well as date of symptoms onset from more than 400 patients. The generalized odds-rate class of regression models has been used to estimate the distribution of the incubation time accounting for age, gender and vaccination type and considering the interval-censored nature of the collected data.

PLENARY LECTURES

TOWARDS RELIABLE EMPIRICAL EVIDENCE IN METHODOLOGICAL BIostatistical RESEARCH: RECENT DEVELOPMENTS AND REMAINING CHALLENGES

Anne-Laure Boulesteix¹

*¹Ludwig Maximilian University of Munich, The Institute for
Medical Information Processing, Biometry, and Epidemiology*

e-mail: boulesteix@ibe.med.uni-muenchen.de

Statisticians are often keen to analyze the statistical aspects of the so-called “replication crisis in science”. They condemn fishing expeditions and publication bias across empirical scientific fields applying statistical methods, such as health sciences. But what about good practice issues in their own - methodological - research, i.e. research considering statistical (or more generally, computational) methods as research objects? When developing and evaluating new statistical methods and data analysis tools, do statisticians and data scientists adhere to the good practice principles they promote in fields which apply statistics and data science? I argue that methodological researchers should make substantial efforts to address what may be called the replication crisis in the context of methodological research in statistics and data science, in particular by trying to avoid bias in comparison studies based on simulated or real data. I discuss topics such as publication bias, cherry-picking, and the design and necessity of neutral comparison studies, and review recent positive developments towards more reliable empirical evidence in the context of methodological biostatistical research.

20 YEARS OF STATISTICAL BIOINFORMATICS: A BRIEF HISTORY OF THE LIMMA AND edgeR PACKAGES

Gordon Smyth¹

¹Walter and Eliza Hall Institute of Medical Research (WEHI)

e-mail: smyth@wehi.edu.au

Genomic technologies underpin modern biomedical research and produce data that is hugely multidimensional. Statistical thinking is highly relevant but classic univariate statistical methods can perform poorly in the high dimensional context. The limma and edgeR software packages were developed for analysing gene expression data from microarray and RNA-seq technologies and have become widely adopted. This talk will give a brief historical perspective on the packages and a discussion of some of the statistical ideas that made them successful.

ACCOUNTING FOR OVERDIAGNOSIS IN ESTIMATING COMPONENTS OF SURVIVAL TIME IN RANDOMIZED CANCER SCREENING TRIALS

Karen Kafadar¹

¹Department of Statistics, University of Virginia

e-mail: *kkuvastat@gmail.com*

Cancer screening is assumed to be beneficial, in terms of reduced mortality and extended survival. Screen-detected survival, when measured as the time between clinical detection of disease and endpoint (cure or death), can be biased by lead time (time by which the screening test advances the time of clinical diagnosis), length biased sampling (probability of detection depends on length), and overdiagnosis (cases that would never have surfaced in the absence of screening). We quantify these effects in this talk and illustrate their non-trivial impacts on the results from actual randomized cancer screening trials. The concepts apply to general periodic screening programs for other conditions and devices.

(This work is performed in collaboration with Dr. Philip C. Prorok, former Chief of the Biometry Research Group, National Cancer Institute.)

SIMULATION APPROACH IN THE DESIGN AND PLANNING CANCER SCREENING TRIALS

Ping Hu¹

¹*Biometry Research Group, National Cancer Institute, USA*

e-mail: pingh@mail.nih.gov

Cancer screening trials aim to evaluate the effectiveness of screening methods in detecting cancer early, improving cure rates, and extending survival. Designing and planning these trials is a complex and challenging task due to the need to balance potential benefits and harms, optimize resource allocation, and account for various uncertainties. The simulation approach offers several advantages over classical probability models, including greater flexibility, better handling of uncertainty, adaptability to new information, optimization capabilities, and opportunities for validation.

The objective of this study is to update and modify the Hu & Zelen model (Biometrika 1997) using simulation methods, making it more flexible and broadly applicable for the design and planning of cancer screening trials. The Hu & Zelen model addresses three fundamental issues in cancer screening trials: 1) the probability of death of an individual from the study or control group, 2) the mortality reduction at each follow-up time, and 3) the power of the statistical test for comparing mortality. This probability model provides outcomes for sample size, mortality reduction, and power.

The updated and modified Hu & Zelen model, developed using simulation methods, includes additional outcomes or features and provides estimations of the number of deaths (mortality) and number of cases (incidence) with or without stratification by stage at each year of follow-up. The validation of the simulation model is applied to the data of the National Lung Screening Trial (NLST). This enhanced model can better inform the design and planning of cancer screening trials, ultimately leading to more effective and efficient screening strategies.

REPLICABILITY ISSUES IN MEDICAL RESEARCH: SCIENCE AND POLITICS

Yoav Benjamini¹

*¹Dept. of Statistics and Operations Research and the Sagol School for Neuroscience
Tel Aviv University*

e-mail: ybenja@gmail.com

Significant testing and the p-value were blamed with replicability problems in medical research. I will argue that selective inference and irrelevant variability are two statistical issues hindering replicability across science. I will review the first in the context of secondary endpoint analysis in clinical and epidemiological research. I will present practical approaches for confidence intervals construction that offer false coverage rate control. These intervals seem to accommodate the concerns of NEJM editors, as reflected in their recent guidelines for authors. I shall discuss addressing the relevant variability in the context of preclinical animal experiments, reporting the results of a large replicability enhancement experiment conducted recently in 3 laboratories.

INVITED TALKS

FUNCTIONAL ADDITIVE MODELS FOR SHAPES AND FORMS OF PLANE CURVES

Almond Stöcker¹, Lisa Steyer², Sonja Greven³

¹*École Polytechnique Fédérale de Lausanne (EPFL)*

^{2,3}*Humboldt-Universität zu Berlin (HU Berlin)*

e-mail: almond.stoecker@epfl.ch

In many imaging data problems, the coordinate system of recorded objects is arbitrary or explicitly not of interest. Statistical shape analysis addresses this by identifying the ultimate object of analysis as the “shape” of an observation, i.e., its equivalence class modulo translation, rotation and re-scaling, or as its “form” modulo translation and rotation. The shape/form space of this equivalence class is endowed with a Riemannian manifold geometry, which needs to be respected in the analysis.

We introduce a flexible additive regression framework for modeling the shape or form of planar (potentially irregularly sampled) curves and/or landmark configurations in dependence on scalar covariates. Models are fit by a novel component-wise Riemannian L2-Boosting algorithm, which yields desirable means of regularization for high-dimensional scenarios and allows estimation of a large number of parameter-intense model terms with inherent model selection.

We utilize the framework A) to analyze configurations of 2D landmarks and outline segments describing the shape of astragali (ankle bones) of wild and domesticated sheep, taking also other “demographic” variables into account, and B) to analyze the form of (irregularly sampled) cell outlines generated from a cellular Potts model in dependence on different metric biophysical model parameter effects (including smooth interactions).

Graphic illustration usually plays an essential role in practical interpretation of smooth (non-linear) additive model effects but becomes a challenging task when the response presents an (equivalence classes of) planar curves or landmark configurations. Therefore, we also suggest a novel visualization for multidimensional functional regression models. Analogous to principal component analysis often used for visualization of functional data, a suitable tensor-product factorization decomposes each covariate effect. After decomposition, main effect directions can be illustrated on the level of curves, while the effect into the respective direction is visualized by standard effect plots for scalar additive models.

ENHANCING CLINICAL TRIAL POWER THROUGH COVARIATE ADJUSTMENT MODELS: AN INVESTIGATION OF METHODS, CLAIMS AND REALISTIC EXPECTATIONS

Aris Perperoglou¹

¹GSK, UK

e-mail: a.perperoglou@gmail.com

Clinical trials are a cornerstone of medical research, driving the development of novel therapies and interventions. They are also still, the most inefficient and ethically debatable aspect of drug development. To ensure successful, safe, and fast delivery of the right treatment to the right patient, there is an increasing demand for efficiency and precision at the design and analysis phase of the trial. In May 2021, the Food and Drug Administration of the USA, published a guidance for industry document, on *Adjusting for Covariates in Randomized Clinical Trials for Drugs and Biological Products*. The paper emphasizes the importance of pre-specifying covariate adjustment methods in the statistical analysis plan to minimize bias and increase the efficiency of clinical trial results, outlines various statistical approaches and highlights the benefits of covariate adjustment in improving precision and reducing sample size. It also discusses potential challenges and considerations in the application of these methods.

Later the same year, in September 2021, the Committee for Medicinal Products for Human Use (CHMP) of the European Medicine Association, qualified a Prognostic Covariate Adjustment (PROCOVA™), a trademarked method that can be regarded as a special case of ANCOVA, as a prognostic score adjustment method that under certain circumstances may increase power or precision of treatment effect estimates in clinical trials with continuous outcomes. With their publication, CHMP did not intend to single out this specific method as ‘the’ method to be used. Instead, they suggest for analysts to compare properties of candidate methods that apply adjustment for covariates and select the one that they consider serving best the objectives of their clinical trial.

In this presentation, we will provide an exploration of the statistical methods and techniques used to create and implement prognostic indices derived from historical data for covariate adjustment in subsequent trials. We will begin by examining the foundations of covariate adjustment models, including various regression techniques, propensity score matching, and machine learning algorithms. We will identify the realistic outcomes that can be achieved under specific circumstances and present simulation results to illustrate the impact and showcase the effectiveness of these methods in improving the power of clinical trials. Finally, we will discuss the future of clinical trials enriched by statistical learning, and how the integration of cutting-edge statistical tools and methods can lead to more efficient, informative, and patient-centered trial designs.

INCORPORATING FUZZY LOGIC PHILOSOPHY IN THE EVALUATION OF FORECASTING MODELS

Cagdas Hakan Aladag¹

¹Department of Statistics, Faculty of Science, Hacettepe University, Türkiye

e-mail: *chaladag@gmail.com*

Making accurate predictions about the future is a challenging task, especially when dealing with uncertainty. In various fields, including economy, health, food, and energy, there are many situations where it's difficult to determine the outcome with certainty. This is because there are numerous variables that can affect the situation, and these variables themselves may not be entirely certain.

To overcome the challenge of modeling under uncertainty, fuzzy logic-based methods have emerged as a viable option. Fuzzy logic allows for the use of membership functions to represent the degree of uncertainty in a particular situation. However, evaluating the performance of fuzzy logic-based forecasting models can be challenging, as they operate differently from classical models.

Most conventional performance measures used to evaluate forecasting models do not take into account membership values, which are a crucial aspect of fuzzy logic-based models. As a result, the evaluation of fuzzy logic-based models using these measures can lead to misleading results. This is because the basic philosophy of fuzzy logic, which incorporates the concept of membership values, is not fully incorporated into the evaluation process.

To address this issue, a performance criterion that takes into account membership values should be used to evaluate fuzzy logic-based models. This approach ensures that the philosophy of fuzzy logic is fully incorporated into the evaluation process, allowing for more accurate predictions in uncertain situations. It's important to note that the evaluation of fuzzy logic-based models should be carried out differently from classical models due to their distinct operating principles.

Overall, incorporating the philosophy of fuzzy logic into the evaluation of forecasting models can lead to more accurate and reliable predictions, especially in situations where there is a high degree of uncertainty. By accounting for the uncertainty in the situation using membership values, fuzzy logic-based models can provide more robust and accurate predictions.

IN SILICO: SIMULATORS AND EMULATORS IN THE EUROPEAN HUMAN BRAIN PROJECT

David M. Steinberg¹, Mira Marcus-Kalish¹, David Refaeli¹, Gilad Shapira¹, Ella Shaposhnik¹

¹*Department of Statistics and Operations Research, Tel Aviv University, Israel*

e-mail: dms@tauex.tau.ac.il

The European Union Human Brain Flagship Project aims to gain an in-depth understanding of the complex structure and function of the human brain with a unique interdisciplinary approach at the interface of neuroscience and technology. It has led to the creation of EBRAINS, a research infrastructure that “facilitates the integration of brain science across disciplines”. One of the major challenges is to model neural activity, from micro- to macro-scale, in a way that enables simulation of the human brain. This leads to so-called *in silico* experiments, which are used “to validate models, and to perform investigations that are not possible in the laboratory”. *In silico* experiments are a particular example of what statisticians describe as *computer experiments*. I will present work from our research team that shows how the related statistical research can be usefully incorporated into the neuroscience context, emphasizing the advantages of replacing a slow/expensive simulator with a much faster and cheaper statistical emulator. This can be especially useful when the simulator runs are matched to data in the context of statistical inference. The talk will present methods for both deterministic simulators, in which the outputs are exactly computed from the inputs, and stochastic simulators, in which each input vector leads to a distribution of outputs. The latter are especially challenging, as they require estimation of the likelihood function. Case studies include the response of neocortical L2/3 large basket cells to electrical stimulation, intracellular reactions stimulated by calcium concentration and several challenging test problems with stochastic simulators.

This research has received funding from the European Union’s Horizon 2020 Framework Programme for Research and Innovation under the Specific Grant Agreement No. 785907 (Human Brain Project SGA2) and the Specific Grant Agreement No. 945539 (Human Brain Project SGA3).

FUNCTIONAL DATA ANALYSES TO ACCOUNT FOR AND INTERPRET GLUCOSE PATTERNS IN CONTINUOUS GLUCOSE MONITORING

Emrah Gecili^{1,2}, Rhonda Szczesniak^{1,2}

¹*Division of Biostatistics & Epidemiology, Cincinnati Children's Hospital Medical Center, OH, USA*

²*Department of Pediatrics, University of Cincinnati, Cincinnati, OH, USA*

e-mail: emrah.gecili@cchmc.org

Identifying phenotypes of type 1 diabetes based on glucose curves from continuous glucose monitoring (CGM) using functional data (FD) analysis can provide valuable insights into the diseases's progression and management. We present a reliable prediction model that can accurately predict glycemic levels using historical data collected from the CGM sensor and real-time risk of hypo-/hyperglycemic for individuals with type 1 diabetes. A longitudinal cohort study was conducted involving 443 type 1 diabetes patients, utilizing CGM data obtained from a completed trial. The FD analysis approach, sparse functional principal components (FPCs) analysis was used to identify phenotypes of type 1 diabetes glycemic variation. We also employed a nonstationary stochastic linear mixed-effects model (LME) that accommodates between-patient and within-patient heterogeneity to predict glycemic levels and real-time risk of hypo-/hyperglycemic by creating specific target functions for these excursions. The first two FPCs explained the majority of the variation in glucose trajectories. The CGM profiles showed higher-order variation during the weeknights, but overall variation was greater on weekends. The employed model has low prediction errors and yields accurate predictions for both glucose levels and real-time risk of glycemic excursions. By identifying these distinct longitudinal patterns as phenotypes, interventions can be targeted to optimize type 1 diabetes management for subgroups at the highest risk for compromised long-term outcomes such as cardiac disease or stroke. Further, the estimated change/variability in an individual's glucose trajectory can be used to establish clinically meaningful and patient-specific thresholds that, when coupled with probabilistic predictive inference, provide a useful medical monitoring tool.

BIostatISTICS AND SARS-COV-2: RESEARCH, POLICY ADVICE, AND COMMUNICATION

Geert Molenberghs¹

¹*Interuniversity Institute for Biostatistics and statistical Bioinformatics
(1) I-BioStat, Hasselt University, Hasselt, Belgium
I-BioStat, KU Leuven, Belgium*

e-mail: Geert.molenberghs@uhasselt.be, geert.molenberghs@kuleuven.be

The COVID-19 pandemic, induced by the SARS-CoV-2 virus, is literally a rare event in the course of history. We need to go back to 1918 for an even worse pandemic, the Spanish Flu, or H1N1, although we also had the tuberculosis pandemic in the interbellum; there was the Russian flu in 1890 (maybe also a coronavirus and *not* influenza), and the plague that literally haunted the world for several centuries.

When there are no antiviral means to speak of, and in the absence of vaccines, time-honored non-pharmaceutical interventions enter stage. Apart from controlling the epidemic, for better or for worse, they generate side effects, for society, its well-being, and for the economy. Based on data and imperfect evidence, the biostatistician contributes to understanding what has happened and is happening, is able to separate signal from noise in predictions for the short- and mid-term future.

Biostatisticians can, and actually should play a role in the response to the COVID-19 crisis, ranging from mathematical and statistical modeling, over day-to-day monitoring, to scientific and government committee work and policy making. We place the mathematical and statistical work done against the background of its use towards policy making, public communication, and outreach in real time. Attention is given to the post-pandemic era, in terms of pandemic preparedness.

ESTIMATING OPTIMAL INDIVIDUALIZED TREATMENT RULES WITH MULTISTATE PROCESSES

Giorgos Bakoyannis¹

¹*Indiana University, USA*

e-mail: gbakogia@iu.edu

Multistate process data are common in studies of chronic diseases such as cancer. These data are ideal for precision medicine purposes as they can be leveraged to improve more refined health outcomes, compared to standard survival outcomes, as well as incorporate patient preferences regarding quantity versus quality of life. However, there are currently no methods for the estimation of optimal individualized treatment rules with such data. In this work, I propose a nonparametric outcome weighted learning approach for this problem in randomized clinical trial settings. The theoretical properties of the proposed methods, including Fisher consistency and asymptotic normality of the estimated expected outcome under the estimated optimal individualized treatment rule, are rigorously established. A consistent closed-form variance estimator is provided and methodology for the calculation of simultaneous confidence intervals is proposed. Simulation studies show that the proposed methodology and inference procedures work well even with small sample sizes and high rates of right censoring. The methodology is illustrated using data from a randomized clinical trial on the treatment of metastatic squamous-cell carcinoma of the head and neck.

RAPIDLY ADVANCING CRISPR SYSTEMS HOLD A GREAT POTENTIAL IN RESEARCH ON DRUG TARGETS

Göknur Giner¹

*¹The Walter and Eliza Hall Institute of Medical Research,
Bioinformatics and Blood Cell and Blood Cancer Divisions*

e-mail: giner.g@wehi.edu.au

In only a few years, as a breakthrough technology, clustered regularly interspaced short palindromic repeats/CRISPR-associated protein (CRISPR/Cas9) gene editing systems have ushered in the era of genome engineering with the plethora of applications, particularly in medicine, biology, pharmacology, and biotechnology.

This talk focuses on some of the most promising CRISPR gene editing tools, such as transcriptional knock-out and activation screens and base editors. Here we discuss how we utilized those technologies in modeling Venetoclax (BCL2 inhibitor) resistance in two different types of cancers and how we handled the computational aspects of analyzing two different types of CRISPR screens from high throughput next-generation sequencing platforms. One of those screens is the CRISPR activation mouse model that is established in our laboratory for inducing gene expression in vivo and in vitro for lymphomas, and the second one is the whole genome knock-out screen to investigate the impact of combination therapies in breast cancer.

Besides recapitulating the tremendous potential held by recent CRISPR technologies in drug development pipelines using aforementioned applications, we also demonstrate our most recent web application that offers a few modules for biologists to explore their own CRISPR data sets to identify and investigate the role of potential drug targets for a given disease of interest.

SPARSE KERNEL FACTOR ANALYSIS MODEL FOR HIGH-DIMENSIONAL UNDERSAMPLED CANCER DATA SETS AND SUPERVISED CLASSIFICATION WITH INFORMATION COMPLEXITY CRITERION

Hamparsum Bozdogan¹

¹*The University of Tennessee, Knoxville, TN 37996 U.S.A.*

e-mail: bozdogan@utk.edu

High-dimensional datasets are common in genomics and medical sciences in cancer detection where the number of features p is much larger than the number of observations n . This is known as the “large p and small n ” problem. Such data sets are often called wide datasets. They pose challenges and a great deal of difficulty to analyze or visualize using standard statistical tools. Most classical statistical methods degenerate in wide dataset scenarios since the covariance matrix or the Gram matrix is ill-conditioned and cannot be inverted to obtain a reliable hyperparameter estimation for statistical inference. To resolve the present existing difficulties, this paper for the first time introduces and presents a novel statistical model, namely Sparse Kernel Factor Analysis (SKFA) in kernel space for dimension reduction and to carry out supervised classification. The SKFA model allows for the analysis of nonlinear features of data by combining the kernel method with the Sparse Factor Analysis. Information complexity (ICOMP) criterion is introduced to learn the hyperparameters of the model simultaneously during the model selection process and dimension reduction. We show our results in the feature space on benchmark wide real-world cancer datasets for dimension reduction and supervised classification. Our analysis confirm the excellent performance of SKFA when the number observations is much smaller than the dimension of the data.

MISSING TIME-DEPENDENT COVARIATE VALUES IN A COX MODEL – JOINT MODELS APPROACH VERSUS COMBINATION OF MULTIPLE IMPUTATION AND JOINT MODELS

Havi Murad¹, Nirit Agay^{1,2}, Rachel Dankner^{2,3}

¹*Biostatistics and Biomathematics Unit, Gertner Institute, Sheba Medical Center, Tel-Hashomer, Israel*

²*Cardiovascular Epidemiology Unit, Gertner Institute, Sheba Medical Center, Tel-Hashomer, Israel*

³*Department of Epidemiology and Preventive Medicine, Sackler Faculty of Medicine, School of Public Health, Tel-Aviv University, Tel-Aviv, Israel*

e-mail: HaviM@gertner.health.gov.il ; Havimurad@gmail.com

We present a novel combination of two approaches used when estimating the association between a time-dependent covariate (marker), measured with missing values, and a survival outcome: multiple imputation (MI) and jointly modelling longitudinal and survival data (JM) (Rizopoulos, 2012).

We have previously developed a procedure for imputing missing values for time-dependent covariates in a discrete time Cox model using the chained equations method (Murad et al., 2020). This time-sequential MI procedure multiply imputes the missing values for each time-period in a time-sequential manner, using completed covariates (imputed and observed) from previous time-periods, but not from future ones, as well as the survival outcome. Recently, we have developed a Fully Conditional Specification MI version that is compatible with the substantive model, in our case the discrete Cox model. Following Bartlett et al. (2015), we term it the Substantive Model Compatible FCS (SMC-FCS). In this method, the missing values are imputed in a chain over all time-periods. In each time-period, the imputation model includes past-imputed values of the marker, as well as future ones in time-periods before the event has occurred, and the survival outcome. Both versions of MI can be applied using the MI procedure in SAS with FCS statement or using similar packages in other software, e.g. the *mice* package in R.

In this presentation, we demonstrate a novel two-step approach, which: (i) multiply imputes the missing values and then (ii) applies JM to each completed data file and combines the estimates using Rubin's rule. In the first step, we present two versions: Time-sequential MI and SMC-FCS MI. We therefore compare three methods: (i) Time-sequential MI + JM (ii) SMC-FCS MI + JM and (iii) a one-stage JM (simple JM). The JM can be executed using the packages *JointModel* or *JointModelBayes* in R.

We use simulations based on data of glucose control variables (plasma glucose and %HbA1c) among diabetic patients, from a large Israeli population-based cohort database (n=546,000) (Dankner et al, 2016), using these methods to evaluate the association of glucose control with the risk of cancer. We examine different patterns of missing data in the glucose control variables (completely missing at random, missing at random and non-missing at random) and the impact of these patterns on the performance of the three methods.

References

- Rizopoulos, D. (2012). *Joint models for longitudinal and time-to-event data: With applications in R*. CRC press.
- Murad, H.**, Dankner, R., Berlin, A., Olmer, L., & Freedman, L. S. (2020). Imputing missing time-dependent covariate values for the discrete time Cox model. *Statistical methods in medical research*, 29(8), 2074-2086. doi: 10.1177/0962280219881168.
- Bartlett, J. W., Seaman, S. R., White, I. R., Carpenter, J. R., & Alzheimer's Disease Neuroimaging Initiative*. (2015). Multiple imputation of covariates by fully conditional specification: accommodating the substantive model. *Statistical methods in medical research*, 24(4), 462-487.
- Dankner, R., Boffetta, P., Balicer, R. D., Keinan-Boker, L., Sadeh, M., Berlin, A., Olmer, L., Goldfracht, M., Freedman, L. S. (2016). Time-dependent risk of cancer after a diabetes diagnosis in a cohort of 2.3 million adults. *American journal of epidemiology*, 183(12), 1098-1106.

STATISTICAL INFERENCE FOR COMPLEX TIME-TO-EVENT DATA UNDER NON-RANDOMIZED COHORTS

KyungMann Kim¹, Tuo Wang¹, Lu Mao¹, Aldo Cocco²

¹*University of Wisconsin-Madison, Madison, WI, USA*

²*Indigo.ai, Milan, Italy*

e-mail: kyungmann.kim@wisc.edu

In multi-season clinical trials with a randomize-once strategy, patients enrolled from previous seasons who stay alive and remain in the study will be treated according to the initial randomization in subsequent seasons. To address the potentially selective attrition from earlier seasons for the non-randomized cohorts, we develop an inverse probability of treatment weighting method to produce unbiased estimates of survival functions or hazard ratios using season-specific propensity scores. Bootstrap variance estimators are used to account for the randomness in the estimated weights and the potential correlations in repeated events within each patient from season to season. Simulation studies show that the weighting procedure and bootstrap variance estimator provide unbiased estimates and valid inferences in Kaplan-Meier estimates and Cox proportional hazard models. Finally, data from the INVESTED trial are analyzed to illustrate the proposed method.

INCORPORATING SHARED PEPTIDES FOR IMPROVED INFERENCE ON THE PROTEINS' ABUNDANCE BASED ON MASS SPECTROMETRY DATA

Malgorzata Bogdan¹

¹*Lund University, University of Wroclaw*

e-mail: *Malgorzata.Bogdan@uwr.edu.pl*

Mass spectrometry (MS) is a core technology for proteomics. It allows the identification and quantification of proteins in biological samples. In a mass spectrometry experiment, peptides – smaller fragments of proteins - are ionized, separated based on their mass and charge, and quantified. Then the relative abundance of proteins is estimated based on the abundance of the respective peptides. This process of the protein quantification becomes difficult in the presence of shared peptides, i.e. the peptides which can belong to several proteins. Typically, shared peptides are removed from analysis of MS data, which leads to loss of peptide-level information and lack of ability to estimate abundances of proteins that are identified only by shared peptides or by a single unique peptide. In this talk we will present a novel method for labeled MS experiments, where multiple biological conditions or subjects may be measured jointly. In this case, peptides have natural quantitative profiles, which we use to estimate the degree of their protein membership (weights). We use these weights to estimate protein-level summaries of peptide data, which are then used for comparisons of the protein abundance under different biological conditions. The results of real data analysis and simulation experiments illustrate the improved precision of the protein quantification, specifically for the proteins who have only few unique peptides. This is a joint research with Mateusz Staniak (University of Wroclaw) and Olga Vitek (Northeastern University).

BRINGING GROUP SPARSITY TO BAYESIAN HIERARCHICAL MODELS: MULTIVARIATE BERNOULLI DISTRIBUTION IS HERE TO HELP

Mehmet Gönen¹

¹*Koç University, Türkiye*

e-mail: *mehmetgonen@ku.edu.tr*

Group sparsity is an important tool to achieve better sparse recovery when there is a group structure in the predictor variables. For example, in bioinformatics applications, we can use our prior knowledge to group the predictor variables and pick the relevant groups that are predictive of the target variable while training the regression or classification algorithm. There are different ways of bringing group sparsity to Bayesian hierarchical models such as the well-known spike-and-slab prior. In this talk, I will introduce our recent and novel solution to this problem using the multivariate Bernoulli distribution under an efficient variational approximation inference scheme. I will mention several extensions of our model and report computational results for different problems, namely, regression, classification, multi-output regression, and multilabel classification.

HOW RELIABLE IS RESTING-STATE FUNCTIONAL CONNECTIVITY?

Xu Meng¹, Philip T. Reiss¹, Ivor Cribben²

¹*University of Haifa*

²*University of Alberta*

e-mail: reiss@stat.haifa.ac.il

Whereas functional magnetic resonance imaging (fMRI) studies examine participants' response to a stimulus or task, resting-state fMRI studies seek to analyze interaction among different brain regions in the absence of an explicit task. Much of resting-state fMRI research focuses on functional connectivity, which may be defined as a simple correlation matrix of brain activity, as measured by the blood-oxygen-level-dependent (BOLD) signal, in a set of brain regions. While such functional connectivity matrices have found widespread use in neuroscience and psychiatry, there is a need for novel methods to assess their reliability, since classical indices of reliability such as the intraclass correlation coefficient (ICC) are designed for scalar- rather than matrix-valued measures. To meet this need, we propose a distance-based ICC (dbICC), defined in terms of arbitrary distances among observations. We introduce a bias correction to improve the coverage of bootstrap confidence intervals for the dbICC, and demonstrate its efficacy via simulation. The Spearman-Brown formula, which shows how more intensive measurement increases reliability, is extended to encompass the dbICC. Our analyses of an existing resting-state fMRI data set suggest that reliability is quite low by conventional standards, irrespective of the choice of distance.

NONPARAMETRIC FUNCTIONAL DATA MODELING OF PHARMACOKINETIC PROCESSES WITH APPLICATIONS IN DYNAMIC PET IMAGING

R. Todd Ogden¹

¹*Columbia University, Mailman School of Public, USA*

e-mail: *to166@cumc.columbia.edu*

Modeling a pharmacokinetic process typically involves solving a system of linear differential equations and estimating the parameters upon which the functions depend. In order for this approach to be valid, it is necessary that a number of fairly strong assumptions hold, assumptions involving various aspects of the kinetic behavior of the substance being studied. In many situations, such models are understood to be simplifications of the "true" kinetic process. While in some circumstances such a simplified model may be a useful (and close) approximation to the truth, in other cases, important aspects of the kinetic behavior cannot be represented. We present a nonparametric approach, based on principles of functional data analysis, to modeling of pharmacokinetic data. We illustrate its use through application to data from a dynamic PET imaging study of the human brain.

REMOVING UNWANTED VARIATION FROM LARGE GENE EXPRESSION DATA WITH RUV-III-PRPS

Ramyar Molania¹

¹*The Walter and Eliza Hall Institute of Medical Research, Bioinformatics Divisions, Melbourne*

e-mail: molania.r@wehi.edu.au

Accurate identification and effective removal of unwanted variation is essential to derive meaningful biological results from RNA sequencing (RNA-seq) data, especially when the data come from large and complex studies. Using RNA-seq data from The Cancer Genome Atlas (TCGA), we examined several sources of unwanted variation and demonstrate here how these can significantly compromise various downstream analyses, including cancer subtype identification, association between gene expression and survival outcomes and gene co-expression analysis. We propose a strategy, called pseudo-replicates of pseudo-samples (PRPS), for deploying our recently developed normalization method, called removing unwanted variation III (RUV-III), to remove the variation caused by library size, tumor purity and batch effects in TCGA RNA-seq data. We illustrate the value of our approach by comparing it to the standard TCGA normalizations on several TCGA RNA-seq datasets. RUV-III with PRPS can be used to integrate and normalize other large transcriptomic datasets coming from multiple laboratories or platforms.

INVESTIGATING DIRECTIONAL CAUSALITY IN MULTICHANNEL BRAIN SIGNALS: THRESHOLD AUTOREGRESSIVE MODELING BASED APPROACH

Sipan Aslan^{1,3}, Hernando Ombao²

¹King Abdullah University of Science and Technology (KAUST), Statistics Program, CEMSE Division

²King Abdullah University of Science and Technology (KAUST), Statistics Program, CEMSE Division

³Van Yuzuncu Yil University (Van YYU), IIBF, Econometrics (On Leave)

e-mail: sipan.aslan@kaust.edu.sa

This study proposes a new analytical outlook to reveal the directional causal relationships in multichannel electroencephalogram (EEG) recordings. EEGs are, in fact, projections of neuronal time-varying phases during any brain activity (e.g., during motor movement/imaging activity). One key goal is to dissecting directionality of such time-dependent transitions. Our approach here is develop a functional form of thresholding through non-linear threshold autoregressive (TAR) models. This leads to identification of potential directional causal relationship as it causes changes in the signal flow. Under the Granger causality framework, causality is mainly identified through predictability of signal. In addition to or beyond this useful G-causality perspective, causality can be further examined by uncovering the mechanism that causes regime changes. In other words, the condition that significantly triggers any sudden (i.e., spikes), temporal, or variable-duration phase changes during an EEG signal can be considered a sign of directional causality. This condition can be identified by estimating the switching mechanism in TAR modeling. From the perspective of this study, it is assumed that the shifting mechanism underlies the form of directional causality. In this talk, causality inferences in multichannel EEGs will be investigated by mainly utilizing the TAR model family. We demonstrate that threshold terms in TAR modeling causing changes between regimes could be considered as interpretable conditions that can define the dynamic directional causality relationship. In addition to presenting exploratory analyses of the publicly available multi-subject Motor Movement/Imagery EEG dataset, we outline the advantages of TAR modeling as a promising method to investigate causality. Furthermore, preliminary results for multivariate TAR analysis are presented in this study.

VALIDATION OF MODEL SELECTION PROCEDURES IN HIGH-DIMENSIONAL ANALYSIS

Victor Kipnis¹, Grant Izmirlian¹, Douglas Midthune¹

¹*Biometry Research Group
USA National Cancer Institute*

e-mail: kipnisv@mail.nih.gov

Model selection is an algorithmic procedure that, within a class of competing models, selects the model by optimizing an estimated model performance criterion. High-dimensional data are particularly prone to overfitting (fitting both signal and noise) leading to two major problems. First, due to presence of noise, a selected model may function well on a training sample but be inaccurate using independent data. Second, due to irregularity of selection procedures, the distribution of the performance criterion is unknown theoretically and could not be consistently estimated by computer-intensive methods. Internal inference becomes impossible necessitating external validation with independent data. Those problems are exemplified considering selection of a classifier for predicting disease using a subset of available biomarkers and using the area under the ROC curve (AUC) as the performance criterion. The conditional distribution of the estimated AUC given training data provides inference for the selected model. This requires only one independent validation sample. Evaluating a model selection procedure itself requires estimation of the unconditional distribution of the estimated AUC or at least its characteristics such as the mean and variance. Therefore, a proper validation of model selection procedures requires multiple independent pairs of training and validation samples. Simulations comparing stepwise and LASSO algorithms within logistic regression and Random Forest exemplify problems with model selection in high-dimensional data when the number of potential biomarkers may exceed the available training sample size and the signal to noise ratio for any single biomarker is relatively low.

SYMPOSIUM TALKS

CONFIDENCE INTERVALS FOR THE WEITZMAN OVERLAPPING COEFFICIENT: THE BINORMAL APPROACH AND ALTERNATIVES

Benjamin Reiser¹

University of Haifa, Faculty of Social Sciences, Department of Statistics

e-mail: reiser@stat.haifa.ac.il

The overlap coefficient (OVL) measures the similarity between two distributions through the overlapping area of their distribution functions. Given its intuitive description and ease of visual representation the development of accurate methods for confidence interval construction can be useful for applied researchers. Inferential procedures that can cover the whole range of distributional scenarios for the two underlying distributions are lacking. Such methods, both parametric and non-parametric are discussed. Parametric approaches based on the binormal model show better performance and are appropriate for use in a wide range of distributional scenarios. Methods are assessed through a large simulation study and are illustrated using a dataset from a study on human immunodeficiency virus-related cognitive function assessment.

THE APPLIED STATISTICAL (DATA) SCIENTIST IN A HIGH- PROFILE AND SOCIETAL ENVIRONMENT - IBS AND EMR -

Geert Molenberghs¹

¹*Interuniversity Institute for Biostatistics and statistical Bioinformatics
(1) I-BioStat, Hasselt University, Hasselt, Belgium
I-BioStat, KU Leuven, Belgium*

e-mail: geert.molenberghs@uhasselt.be

A perspective will be offered on the profession of the biometrician, the biostatistician, and more generally the applied statistical scientist, in an ever changing environment. The specifics of working in a multi-disciplinary environment will be discussed, referring to collaboration with agronomists, biologists, epidemiologists, medical professionals, etc. At the same time, interactions with other semi- or fully quantitative fields will be touched upon, such as computational biologists, computer scientists, engineers, etc. The current-day (r)evolution towards data science will be placed against a historical timeline of our field, which saw, over a relatively brief period of just one century, the coming of epidemiology and observational studies, (statistical) genetics, bioinformatics, the omics, big data, data science, data analytics, etc. Historical notes related to the International Biometric Society and its Regions, especially the Eastern Mediterranean Region, will be interwoven.

Reference

Molenberghs, G. (2005). Presidential Address: XXII International Biometric Conference, Cairns, Australia, July 2004: Biometry, biometrics, biostatistics, bioinformatics. *Bio-X. Biometrics*, **61**, 1-9.

CARDIAC CLASSIFICATION USING MACHINE LEARNING MODELS WITH INFORMATION COMPLEXITY

Hamparsum Bozdogan¹

¹*The University of Tennessee, Knoxville, TN 37996 U.S.A.*

e-mail: bozdogan@utk.edu

High-dimensional data is prevalent in many application areas such as in genomics, medical diagnosis, imaging, cancer detection, cyber security, and transportation to mention a few. This paper will present the performance of the select machine and deep learning neural network unsupervised and supervised algorithms using the information complexity (ICOMP) criterion. State-of-the-art machine and deep learning neural network classifiers are introduced to classify risky heart attack patients in medicine in early detection of the cause of heart attack from nuclear magnetic resonance (NMR) images. Among these classifiers, an adaptive kernel Gaussian mixture-model (AK-GMM) cluster analysis; hybrid radial basis function neural networks (HRBF-NN) classifiers; robust Bayesian sparse relevance machines (RBS-RVMs) classifiers; and deep learning neural network (DL-NN) autoencoder classifiers will be presented and their performance will be compared. We show that our approach improves the probability of misclassification error over the existing conventional methods. Our approach provides an expert model-based data-driven approach to healthcare analytics applications in complex pattern recognition problems in many cross-disciplinary fields.

MAPPING THE DISEASE PREVALENCE OF TURKIYE

Mehmet Koçak¹

¹*Istanbul Medipol University*

e-mail: mehmetkocak@medipol.edu.tr

High-dimensional datasets are common in genomics and medical sciences in cancer detection where the number of features p is much larger than the number of observations n . This is known as the “large p and small n ” problem. Such data sets are often called wide datasets. They pose challenges and a great deal of difficulty to analyze or visualize using standard statistical tools. Most classical statistical methods degenerate in wide dataset scenarios since the covariance matrix or the Gram matrix is ill-conditioned and cannot be inverted to obtain a reliable hyperparameter estimation for statistical inference. To resolve the present existing difficulties, this paper for the first time introduces and presents a novel statistical model, namely Sparse Kernel Factor Analysis (SKFA) in kernel space for dimension reduction and to carry out supervised classification. The SKFA model allows for the analysis of nonlinear features of data by combining the kernel method with the Sparse Factor Analysis. Information complexity (ICOMP) criterion is introduced to learn the hyperparameters of the model simultaneously during the model selection process and dimension reduction. We show our results in the feature space on benchmark wide real-world cancer datasets for dimension reduction and supervised classification. Our analysis confirm the excellent performance of SKFA when the number observations is much smaller than the dimension of the data.

GENESELECTML: A COMPREHENSIVE WAY OF GENE SELECTION FOR RNA-SEQ DATA VIA MACHINE LEARNING ALGORITHMS

Ozlem Ilk¹

¹*Department of Statistics, Middle East Technical University, Türkiye*

e-mail: oilk@metu.edu.tr

Selection of differentially expressed genes (DEGs) is a vital process to discover the causes of diseases. It has been shown that modelling of genomics data by considering relation among genes increases the predictive performance of methods compared to univariate analysis. However, there exist serious differences among most studies analyzing the same dataset for the reasons arising from the methods. Therefore, there is a strong need for easily accessible, user-friendly, and interactive tool to perform gene selection for RNA-seq data via machine learning algorithms simultaneously not to miss DEGs. We develop an open-source and freely available web-based tool for gene selection via machine learning algorithms that can deal with high performance computation. This tool includes six machine learning algorithms having different aspects. Moreover, the tool involves classical pre-processing steps; filtering, normalization, transformation, and univariate analysis. It also offers well-arranged graphical approaches; network plot, heatmap, venn diagram, and box-and-whisker plot. Gene ontology analysis is provided for both mRNA and miRNA DEGs. The implementation is carried out on Alzheimer RNA-seq data to demonstrate the use of this web-based tool. Eleven genes are suggested by at least two out of six methods. One of these genes, hsa-miR-148a-3p, might be considered as a new biomarker for Alzheimer's disease diagnosis. GeneSelectML is distinguished in that it simultaneously uses different machine learning algorithms for gene selection and can perform pre-processing, graphical representation, and gene ontology analyses on the same tool. This tool is freely available at www.softmed.hacettepe.edu.tr/GeneSelectML.

* This is a joint work with Osman Dag, Merve Kasikci, and Metin Yesiltepe.

VARIATIONAL MULTIPLE IMPUTATION IN HIGH-DIMENSIONAL REGRESSION MODELS WITH MISSING RESPONSES

Recai M. Yücel¹

¹Temple University, Department of Epidemiology and Biostatistics, College of Public Health

e-mail: recai.yucel@temple.edu

Multiple imputation has become one of the standard methods in drawing inferences in many incomplete data applications. Applications of multiple imputation in relatively more complex settings such as high-dimensional clustered data require specialized methods to overcome the computational burden. Using mixed-effects models, we develop methods that can be applied to continuous, binary, or categorical incomplete data. We overcome the computational burden by employing variational Bayesian inference for sampling the posterior predictive distribution missing data. These methods specifically target high-dimensional covariates and work with spike-and-slab priors, which force the variables of importance to be in the imputation model. The individual regression computation is then incorporated in an increasingly popular variable-by-variable imputation algorithm. Finally, we use calibration-based algorithms to adopt these methods to multiply-impute categorical variables. We present a simulation study to assess the performance of these methods in a repetitive sampling framework.

ERGUN KARAAĞAOĞLU & TURKISH JOURNAL OF BIOCHEMISTRY

Yahya Laleli¹

¹*Duzen Laboratories Group, Türkiye*

e-mail: ylaleli@duzen.com.tr

Dear Ergun, the training and guidance you have provided to the biochemistry community, including myself, will not be forgotten.

I don't even remember when our relationship with you in the Turkish Biochemistry Journal Editorial Board has started, it should be at least 15 years! We used to work in such harmony ... You did contribute your statistical approaches and practices with your suggestions so that the decisions at the board added even higher value to the primary effective output of the articles! Although our journal was the Journal of Medical Biochemistry, we were also accepting articles on industrial production and cleansing of wastes using basic biochemistry and biochemical pathways as research articles in those days. In other words, the articles that reached were examined within the framework of the concepts of both medical statistics and biostatistics under the responsibility of our statistics editor, Ergun Hoca. Test awareness, diagnosis or screening test, evaluation of bias, selectivity, negative and positive likelihood ratios entered our routine concept. Besides quantitative changes (\bar{x} and SD); categorical changes, using frequency and proportion, visual expression of quantitative changes, primary outcome variability, controlling the normal distribution for explanatory variability, making the comparison over median values for the data which did not show normal distribution, searching for 95% confidence interval to determine compatibility between parameters. These were all applied in your leadership together with your colleagues at your department.

I remember very well that apart from giving guiding suggestions for the structuring of the publications within the framework of these views, you suggested "send your raw data and we will evaluate it" for many articles, its constructiveness, the main purpose of teaching, it was. Of course, because complementary data could not be obtained, the rejection was made after the fourth evaluation! What did that mean? It is easy to give reasons and reject, but for you the goal was suggesting an approach is to show direction. You continued this effort without giving up.

As for your academic life, it is not for me to measure your success in education, the publications/products of your graduates are obvious. Also, a concept I have witnessed is that Ergun Hoca's graduates are closely related to the effective continuity of the cultural and social assets as they are needed in research.

How this all came to an end is another tale. When the new Editorial Board demanded that I resign from my editorial position, he said "editing is a group cohesion job" and included himself in the group of 7 who accepted me, not sided on the four who left. It is a concept that needs to be measured how this situation has an impact on article writing trainings for researchers who will start their new publication life.

I am sure he still has a lot to offer. His vision and activities will shed light on us and guide us.

ORAL PRESENTATIONS

OP1. ESTIMATION OF AVERAGE CAUSAL EFFECT IN CLUSTERED DATA WITH COVARIATE MEASUREMENT ERROR

Recai M. Yuçel¹, Raina E. Josberger², Meng Wu³

¹Temple University, Department of Epidemiology and Biostatistics, College of Public Health

²New York State Department of Health

³New York State Department of Health

e-mail: recai.yucel@temple.edu

Statistical inferences using potentially rich and inexpensive administrative data require caution as it is not collected/maintained for inferential purposes. One of the major challenges of administrative data relates to measurement issues as well as generalizability. In this paper, we consider the problem of measurement error which often arises from inaccurate data observations (e.g., self-administered questionnaires) or measurement process such as recall bias. In particular, we aim to conduct causal inference when measurement errors are confined to covariates. We gauge the measurement error using validation data and work with regression calibration to conduct inferences. While regression calibration is one of the most frequently used method for association studies with error-prone covariates, it is rarely applied in causal effect studies. In this paper, we extend regression calibration to draw causal inference in clustered data. The estimation of average causal effect (ACE) is constructed under the potential-outcomes framework and implemented by the model-based approach modified under regression calibration. To overcome analytical and theoretical challenges in the computation of the estimate of ACE variance, we use bootstrap methods for clustered data. We apply these methods to study the ACE of prenatal care on birth weight for lower income women who were insured by New York's Medicaid program.

Keywords: Causal inference, Measurement error, Regression calibration, Low birth weight, Prenatal care

OP2. VIRAL LOAD DYNAMICS OF SARS-COV-2 DELTA AND OMICRON VARIANTS FOLLOWING MULTIPLE VACCINE DOSES AND PREVIOUS INFECTION

Naama M. Kopelman¹, Yonatan Woodbridge^{1,2}, Sharon Amit³, Amit Huppert^{2,4}

¹*Department of Computer Science, Holon Institute of Technology, Holon, Israel*

²*The Gertner Institute for Epidemiology & Health Policy Research, Sheba Medical Center, Ramat Gan, Israel*

³*Clinical Microbiology, Sheba Medical Center, Ramat Gan, Israel*

⁴*Sackler Faculty of Medicine, Tel-Aviv University, Tel-Aviv, Israel*

e-mail: naamako@hit.ac.il

An important aspect of vaccine effectiveness is its impact on pathogen transmissibility, harboring major implications for public health policies. As viral load is a prominent factor affecting infectivity, its laboratory surrogate, qRT-PCR cycle threshold (Ct), can be used to investigate the infectivity-related component of vaccine effectiveness. While vaccine waning has previously been observed for viral load during the Delta wave, less is known regarding how Omicron viral load is affected by vaccination status, and whether vaccine-derived and natural infection protection are sustained. By analyzing results of more than 460,000 individuals, we show that while recent vaccination reduces Omicron viral load, its effect wanes rapidly. In contrast, a significantly slower waning rate is demonstrated for recovered COVID-19 individuals. Thus, while the vaccine is effective in decreasing morbidity and mortality, its relatively small effect on transmissibility of Omicron (as measured here by Ct) and its rapid waning call for reassessment of future booster campaigns.

Keywords: SARS-CoV-2, Viral vaccines, Viral load, Vaccine effectiveness, Regression analysis

OP3. EVALUATING UNIVARIATE, MULTIVARIATE REFERENCE INTERVAL METHODS: A COMPARATIVE ANALYSIS

Esra Kutsal Mergen¹, Sevilay Karahan²

¹ *Hacettepe University Faculty of Medicine, Department of Biostatistics, Ankara, Türkiye*

e-mail: esrakutsalkaynar@gmail.com

Reference intervals are crucial in healthcare for interpreting laboratory test results and making medical decisions. However, univariate reference intervals ignore multivariate relationships, resulting in decreased accuracy and an increased probability of false positives. For instance, a healthy individual has a 95% chance of being classified as healthy with univariate intervals, but only a 90.25% chance with two interrelated variables. To overcome these issues, this study proposes two multivariate reference interval approaches that consider the relationship between variables, reduce the probability of false positives, and alleviate the uncertainty of interpreting the concept of the multivariate reference region. A simulation study was conducted in RStudio using different scenarios for the interrelated Ferritin and Transferrin Saturation and hemoglobin variables used in the diagnosis of iron deficiency anemia to demonstrate the effectiveness of the proposed approaches. Although the multivariate reference region method has been recommended for years, it is still complicated to apply in clinical practice due to difficulties in obtaining and analyzing data. Therefore, this study proposes more accessible multivariate reference interval approaches that reduce the false positive rate and are easy to apply in the clinic. The simulation study showed that the Mahalanobis approach is recommended when facing uncertainty in the multivariate reference region and the difficulty of implementation, as it produced a patient classification rate closest to 5%. However, presenting laboratory results to physicians is still a significant challenge in the multivariate reference interval approach. While there are many unexplored methods, the proposed multivariate reference interval approaches offer a promising area for further research, encouraging clinical laboratory staff and physicians to use them. While Mahalanobis provides the closest patient classification rate to 5%, the Multivariate confidence interval approaches yields the most comparable results and is easier for physicians to evaluate since it does not provide lower and upper reference limits.

Keywords: Mahalanobis distance, Reference interval, Multivariate reference region, Univariate reference interval, Simulation study

OP4. EVALUATION OF OBJECTIVE STRUCTURED EXAMINATION TOOL WITH CLASSICAL TESTING INSTITUTION, GENERALIZABILITY THEORY AND ITEM RESPONSE THEORY

M.Yasemin Akşehirli Seyfeli¹, Atilla H. Elhan², Zeynep Baykan³,
Gözde Ertürk Zararsız⁴, Orhun Öztürk⁵, Gökmen Zararsız⁴, Ahmet Öztürk⁴

¹Erciyes University Medical Faculty Department of Biostatistics

²Ankara University Medical Faculty Department of Biostatistics

³Erciyes University Medical Faculty Department of Medical Education

⁴Erciyes University Medical Faculty Department of Biostatistics

⁵Hacettepe University Department of Statistics

e-mail: yaseminseyfeli@gmail.com

The aim of this study is to compare the validity and reliability analyzes of classical test theory, generalizability theory and item response theory by using the data of a 24-item assessment guide used in the osce exams at the Faculty of Medicine, and to discuss their advantages and disadvantages. In the study, an exam data in the first term 340 students, and, the second term 328 students in total by filling in the 24-item assessment guide as they did/did not complete each item was used. According to the results of the explanatory factor analysis, 50 percent of the variance sources are explained with the 3-factor structure. With the confirmatory factor analysis, the model fit goodness of these three factor structures were found to be in acceptable ranges ($X^2/sd=2.006$; $RMSEA =0,055$; $CFI =0.942$). When the structure was evaluated with generalizability analysis, it was found that the number of items explained 7.1 percent and 22.4 percent of the individual variable. The G coefficient of the model was estimated as 0.88 and the Phi coefficient as 0.87. In item response theory, it was found that the guide was in a 3-dimensional structure, which did not provide the unidimensionality assumption. The items of each dimension of these three dimensions were examined with the Rasch model and the 2PL model. The items of the first and third factors did not fit the PC analysis. The 10-item model of the second factor conformed to the 2PL model. The empirical reliability coefficients and fit coefficients of the models are lower than the coefficients of the classical test theory. When the evaluation guide, which was found to be valid and reliable, was examined with generalizability, since the evaluation is made with the total correlations of the items in the classical test theory models, it has been determined that the student scores are affected by factors other than the item and student ability. When we look at the item response theory models, the incompatibility of two factors and the fact that one factor only complies with the 2PL model leads us to the conclusion that the validity and reliability of the assessment guide should be improved.

Keywords: IRT, RASCH, G- teory

OP5. REPLICABILITY ACROSS MULTIPLE STUDIES

Ruth Heller¹, Marina Bogomolov²

¹*Department of Statistics and Operations Research, Tel-Aviv University*

²*Faculty of Data and Decision Sciences, Technion*

e-mail: ruheller@gmail.com

Meta-analysis is routinely performed in many scientific disciplines. This analysis is attractive since discoveries are possible even when all the individual studies are underpowered. However, the meta-analytic discoveries may be entirely driven by a single study, and thus non-replicable. The lack of replicability of scientific findings has been of great concern in the last two decades. In order to alleviate the concern, we suggest that a standard meta-analysis shall be complemented with an analysis towards replicability of findings. We address the common setting in modern applications, in which multiple hypotheses are examined in each study. We start by formally defining replicability, and the directional no-replicability error measures we want to control. We suggest the cross-screening approach, and provide procedures that control the directional error measures of interest. This approach uses each study in order to plan the multiple testing procedure of the other study (or studies), as well as in order to test the hypotheses in that study. We demonstrate the usefulness of this approach in a high-dimensional genomic application. We also discuss some of the current challenges.

The relevant paper is [arXiv:2210.00522](https://arxiv.org/abs/2210.00522).

Keywords: High-dimensional studies, Meta-analysis, Multiple comparisons, Replicability

OP6. MODELLING LONGITUDINAL COGNITIVE TEST DATA WITH CEILING EFFECTS AND LEFT SKEWNESS

Denitsa Grigorova¹, Dean Palejev², Ralitz Gueorguieva³

¹*Big Data for Smart Society Institute, Sofia University 125 Tsarigradsko Shosse, Bl. 2, 1113 Sofia, Bulgaria and*

Faculty of Mathematics and Informatics, Sofia University 5 James Bourchier Blvd., 1164 Sofia, Bulgaria

²*Institute of Mathematics and Informatics, Bulgarian Academy of Sciences Acad. G. Bonchev St., Bl. 8, 1113 Sofia, Bulgaria*

³*Department of Biostatistics, Yale School of Public Health, 60 College St, New Haven, CT 06520, U.S.A*

e-mail: denitsa.grigorova@gate-ai.eu, dgrigorova@fmi.uni-sofia.bg

Cognitive tests are among the markers for the development of cognitive diseases such as Alzheimer's disease. We model the scores from the Mini Mental State Examination (MMSE, discrete values in the range 0-30) over time on data from the Alzheimer's Disease Neuroimaging Initiative (ADNI, <http://adni.loni.usc.edu/>). The challenge of modelling such an outcome as MMSE is that the data are left-skewed with ceiling effect - the maximum possible score on the MMSE is 30 and this maximum is often achieved by healthy individuals (the higher the value of MMSE the better cognitive function of the individual). Different approaches for modeling MMSE have been considered in the statistical literature, such as linear mixed effects models on transformed data, mixture models based on latent class growth analysis and generalized additive models for location, scale and shape (GAMLSS). We find models such as binomial and beta-binomial from GAMLSS more appropriate for the MMSE score and we apply them to the ADNI data. We use random effects (parametric and non-parametric) in the models to account for correlations among repeated measures on the same individual. The estimation is based on the maximum likelihood approach. Using Bayesian Information Criterion, we select the best model that fits the data. Additionally, we propose a bootstrap approach for estimation of the covariance matrix of the estimates. Using statistical tests and inference we compare the cognitive function over time for individuals with cognitive impairment and normal controls in the ADNI data set. We also perform simulation studies with different sample sizes that evaluate the binomial and beta-binomial models in terms of bias and efficiency.

Keywords: Alzheimer's Disease Neuroimaging Initiative (ADNI), GAMLSS models, Mini Mental State Examination (MMSE), Random effects

OP7. BRANCHING MODELLING OF MUTATIONS AND RISK ASSESSMENT IN CANCER RESEARCH

Maroussia Slavtchova-Bojkova¹, Kaloyan Vitanov²

¹*Sofia University "St. Kl. Ohridski", Faculty of Mathematics and Informatics and Institute of Mathematics and Informatics at Bulgarian Academy of Sciences*

²*Sofia University "St. Kl. Ohridski", Faculty of Mathematics and Informatics*

e-mail: Bojkova@fmi.uni-sofia.bg

The multi-type decomposable Sevastyanov branching process (MDSBP) represents a model in which two classes of cell types - class W_0 (we usually assume that at least one cell type from W_0 is supercritical) whose cell types can only have offspring that is again with types within W_0 and class W_e with cell types (usually assumed subcritical) that can have offspring from W_e but also can "emit" mutant daughter cells towards types from W_0 . Our motivation for considering this particular setting is that it explores an irreversible path in the evolution of the population since biological populations usually do not revert to less adapted states if conditions remain unchanged. A system of integral equations for the probability generating functions of the process are obtained and accordingly the probabilities of extinction, number of occurred mutations, waiting time to escape mutant, and immediate risk of escaping extinction are studied. We also provide a general numerical scheme for calculating obtained systems of integral equations. Our theoretical results do not depend on particular assumptions regarding lifespan distributions or reproduction rates for individual types, nor on an assumption that the probabilities of mutation are small. These virtues, together with the fact that the MDSBP is in continuous time, distinguish the MDSBP from the models studied before. Furthermore, the model allows dependence from individual cell age. The MDSBP is readily applicable to many biological contexts such as cancer evolution and treatment in the presence of multiple metastases, viruses developing resistance to vaccines, migration of individuals, etc.

Keywords: Decomposable branching processes, Continuous time, Mutations

OP8. JOINT SPATIOTEMPORAL MODELLING OF HUMAN IMMUNODEFICIENCY VIRUS AND TUBERCULOSIS IN ETHIOPIA USING A BAYESIAN HIERARCHICAL APPROACH

Legesse Kassa Debusho¹, Leta Lencha Gemechu²

¹*Department of Statistics, College of Science, Engineering and Technology, University of South Africa, South Africa*

²*Department of Statistics, College of Natural and Computational Sciences, Dire Dawa University, Ethiopia*

e-mail: debuslk@unisa.ac.za

Understanding of the joint geographical patterns of human immunodeficiency virus (HIV) and tuberculosis (TB) infections over time is essential since it helps to target high-risk areas with effective control measures. Therefore, the main objective of the study was to assess the spatial and temporal patterns of HIV and TB burden in Ethiopia jointly at district level for a four-year period, 2015 to 2018. The Bayesian hierarchical joint spatiotemporal modelling was applied to analyse the data. Six models with different priors for the precision of random effects variances were fitted and best fitting model was selected using the DIC. The annual TB case notifications in Ethiopia had inconsistent trend but the aggregated annual number of HIV patients enrolled in HIV care was increased in the first three years and decreased in the last year of the study period. The results from selected model show that about 53% of the variability in HIV and TB incidences in the study period was explained by the shared temporal component, disease-specific spatial effect of HIV, and space-time interaction effect. The shared temporal trend and disease-specific temporal trend of HIV risk had a slight upward trend between 2015 and 2017 then a slight decrease in 2018. However, the disease-specific temporal trend of TB risk had almost constant trend with minimal variation over the study period. The distribution of the shared relative risks was similar with the distribution of disease-specific TB relative risk whereas that of HIV had more districts as high-risk areas. The findings could provide valuable information to health policymakers in Ethiopia in the improvement of the national or district responses for geographically targeted and integrated interventions to jointly control the two diseases.

Keywords: Bayesian hierarchical model, HIV, Poisson regression, Joint spatiotemporal modelling, Tuberculosis

OP9. COMPARING FREQUENTIST AND BAYESIAN APPROACHES FOR MIXED DESIGN ANOVA IN REPEATED MEASUREMENTS: A SIMULATION STUDY WITH EXPONENTIAL DISTRIBUTIONS

Zeynep Özel¹, Ebru Kaya Başar², Mustafa Agah Tekindal¹

¹*Department of Biostatistics, Izmir Katip Çelebi Üniversitesi, Faculty of Medicine, İzmir*

²*Lecture Dr., Akdeniz University, Statistics Consultancy Application and Research Center, Antalya*

e-mail: zozel4225@gmail.com

Mixed design analysis of variance in repeated measurements allows the examination of data obtained from the same measurement performed multiple times. Bayesian and frequentist approaches are two different methods used for estimating a parameter for repeated measurements. In this study, the results obtained from bayesian and frequentist approaches for repeated measurements in balanced and unbalanced samples were investigated. The aim of this study is to compare the results of both frequentist and bayesian approaches using the mixed design ANOVA method for data obtained from four repeated exponential distributions with different coefficient of variations and balanced/unbalanced designs (4x4). Datasets with coefficient of variation of 0.5 and 1 were generated from the exponential distribution for the study design, which consists of balanced and unbalanced samples. The inclusion of unbalanced design aims to achieve positive heterogeneity by increasing the variance. The simulation study was repeated 1000 times. Evaluations were made using both frequentist and bayesian approaches, and the results were evaluated according to confidence intervals, and the approaches were compared. Bayesian and frequentist approaches offer different methods for estimating a parameter in repeated measurements. Both methods have advantages and disadvantages. The choice of appropriate method may depend on the availability of prior information; Bayesian approach may yield more accurate results if prior information is available, while frequentist approach may be more appropriate if no prior information is available. In conclusion, both approaches can be used as alternatives to each other in repeated measurements.

Keywords: Mixed Design ANOVA, Bayesian, Frequentist, Simulation

OP10. DATA-DRIVEN SIMULATIONS FOR QUANTITATIVE BIAS ANALYSES IN REAL-WORLD SURVIVAL ANALYSES

Michal Abrahamowicz¹, Marie-Eve Beauchamp¹, Anne-Laure Boulesteix²,
Tim Morris³, Willi Sauerbrei⁴, Jay Kaufman¹

¹*McGill University, Montreal, Canada*

²*LMU Munich, Munich, Germany*

³*MRC Clinical Trials Unit at UCL, UCL, UK*

⁴*Medical Center - University of Freiburg, Freiburg, Germany*

e-mail: michal.abrahamowicz@mcgill.ca

We propose a novel approach, based on data-driven simulations, to assess the impact of a particular analytical limitation of a specific real-world prognostic or epidemiological study on its results and conclusions. We first describe the 7 steps necessary to implement our approach. Then, we focus on two challenges often encountered in real-world time-to-event (survival) analyses. For each example, we illustrate how data-driven simulations, designed to accurately reflect the salient characteristics of the corresponding real-world dataset, can yield new insights. To this end, we use the permutational algorithm, validated for simulating event times conditional on time-varying exposures and/or effects, to preserve both the empirical distribution of event times and the assumed or observed hazard ratios for all relevant variables. The first illustration concerns assessing the impact of omitting an important prognostic factor (cancer stage) on the results of multivariable Cox proportional hazards (PH) analyses used to explore the independent association of the exposure (colon obstruction) with the hazard of all-cause mortality, among persons diagnosed with colon cancer. Here, we show how data-driven simulations permit assessing the joint impact of (i) unmeasured confounding bias and (ii) non-collapsibility, while separating their effects. The second illustration focuses on the pharmaco-epidemiological study of the association between recent use of sedative medications, modeled as a time-varying exposure, and the incidence of cognitive dysfunction. The event is interval-censored as it can be detected only at discrete times of medical visits. Data-driven simulation results indicate how the strength of a systematic bias toward the null varies depending on the strategy for imputing (unknown) exact event times and on the assumed strength of the underlying association. Our methods and results extend the simulation-based Quantitative Bias Analyses to multivariable time-to-event analyses with time-varying exposures.

Keywords: Survival analysis, Simulations, Bias, Time-varying covariates

OP11. CONDITIONAL RANDOMIZATION TEST FOR AVERAGE TREATMENT EFFECT WITH SURVIVAL FOREST

Mehmet Ali Kaygusuz¹, Vilda Purutçuoğlu

¹*Anadolu University, Department of Economics, 26780, Türkiye*

²*Middle East Technical University, Department of Statistics, 06800, Türkiye*

e-mail: makaygusuz1988@gmail.com

Heterogeneous treatment effect has a fundamental role and a broad range of applications in genetics, epidemiology, and econometrics. Nevertheless, this type of applications can include censored survival outcomes. In this situation, heterogeneous effect may not capture hazard or survival functions. Subsequently, in this study, we consider survival analysis which can be very crucial for estimation with heterogeneous treatment effect on biological networks. But, when we work with censored data, in particular, under high number of parameter (p) with respect to number of observation (n), the estimation of heterogeneous effect with survival outcomes can be computationally difficult. Thus, we use causal survival forest for biological networks to gain from computational efficiency. Furthermore, we apply conditional randomization test since it is known from previous studies that it is promising to define relationship on heterogeneous treatment effect with survival outcomes. We examine proposed model selection procedure to show its efficiency with the different p and n .

Keywords: Multiple testing, Random forest algorithm, Causal models, Survival analysis

OP12. RECONSTRUCTING SURVIVAL DATA FROM PUBLISHED KAPLAN-MEIER CURVES

Georgia Rompoti^{1,2}, Dimitris Karlis³, Urania Dafni^{1,2}

¹National and Kapodistrian University of Athens, Athens, Greece

²Frontier Science Foundation-Hellas, Athens, Greece

³Athens University of Economics and Business, Athens, Greece

e-mail: grompoti@frontier-science.gr

The necessity of reliable and valid decision-making regarding public health issues rendered the implementation of meta-analyses imperative. In view of that, the reconstruction of individual survival data from Kaplan-Meier (KM) curves is highly important. After a literature review, the most robust method (Guyot et al., 2012) was chosen and optimized. Subsequently, a simulation of 50 replications, based on the Weibull distribution, was carried out. The quality of the reconstructed data was examined under various alternative scenarios concerning both the number of points extracted from the initial KM graphs and the number of time points at which the size of the risk set is known. Moreover, the consistency between the initial summary measures and the ones occurring from the reconstructed data was investigated. The results have shown that the larger the number of the extracted coordinates from the initial KM curves is, the closer the algorithm estimate from the true value is. The discrepancies among the estimates when the number of time points (at which the number of subjects at risk is reported) changes have been negligible. However, it was evident from the KM curves that the termination of a study of interest occurred earlier in the reconstructed data. Finally, the hazard ratio estimates in the applied scenarios were consistent and robust, especially when additional information was provided. On the contrary, Restricted Mean Survival Time differences estimates exhibited great variability. Consequently, both the early termination in the reconstructed data and the poor estimates of some of the summary measures require further investigation and careful consideration. Additionally, an in-depth examination of the aforementioned method's accuracy is of the utmost importance on occasions where different distributions are used and under the presence of censoring, two scenarios that reflect better real-life problems.

Keywords: Kaplan-Meier, reconstruction, HR, RMST

References

- Guyot, P., Ades, A. E., Ouwens, M. J., & Welton, N. J. (2012). Enhanced secondary analysis of survival data: reconstructing the data from published Kaplan-Meier survival curves. *BMC medical research methodology*, 12, 1-13.
- Liu, N., Zhou, Y. & Lee, J. (2021). IPDfromKM: Reconstruct Individual Patient Data from Published Kaplan-Meier Survival Curves. *BMC Medical Research Methodology*.

OP13. METHODOLOGICAL ISSUES WITH PROPORTIONAL HAZARD MODELS

Maria-Tereza Dellaporta^{1,2}, Dimitris Karlis¹, Urania Dafni^{2,3}

¹*Athens University of Economics and Business, Athens, Greece*

²*Frontier Science Foundation-Hellas, Athens, Greece*

³*National and Kapodistrian University of Athens, Athens, Greece*

e-mail: mtdellaporta@frontier-science.gr

In recent years, non-proportional data are frequently encountered. Clinical trial data may deviate from the usual assumption of a constant hazard ratio over time as a result of the administration of novel medicinal products and the implementation of innovative therapeutic procedures with unprecedented mechanisms of action. A big part of survival analysis is mainly based on two well-known methods: the log-rank test for the comparison of survival curves and the Cox proportional hazard (PH) model for the estimation of the effect corresponding to numerous variables of interest. Both methods are based on the assumption of proportional hazards and thus, when it is violated they are expected to perform poorly and yield biased results. In this work, various tests for the proportional hazards assumption and numerous testing procedures for the significance of treatment effect, are being compared in terms of performance. Two simulation studies, one for each group of tests, are conducted with five scenarios each: one scenario under proportionality and four with different patterns of non-proportional hazard, usually reported in contemporary publications (early and late effect, crossing hazards, long-term survivors). We show that among eighteen tests for proportionality, Grambsch & Therneau's (1994), as well as the suggestions in Lin (1991), exhibit the best overall performance. Moreover, the comparison of twenty tests for treatment effect illustrates the superiority and flexibility of the Cauchy combination of change-point Cox regressions (Zhang et al., 2021), a newly developed method which also provides piecewise hazard ratio estimates. Nevertheless, none of the tests surpasses all the others under the different alternative hypotheses assumed and therefore, further research is needed expanding the array of tests and the non-PH scenarios.

Keywords: Proportionality, Treatment effect, Test

References

- Grambsch, P. M. and T. M. Therneau (1994). Proportional hazards tests and diagnostics based on weighted residuals. *Biometrika*, 81 (3), 515-52.
- Lin, D. (1991). Goodness-of-fit analysis for the Cox regression model based on a class of parameter estimators. *Journal of the American Statistical Association*, 86 (415), 725-728.
- Zhang, H., Q. Li, D. V. Mehrotra, and J. Shen (2021). CauchyCP: A powerful test under non-proportional hazards using Cauchy combination of change-point Cox regressions. *Statistical Methods in Medical Research*, 30 (11), 2447–2458.

OP14. STACKING BASED APPROACHES FOR SURVIVAL ANALYSIS OF RNA-SEQUENCING DATA

Ahu Cephe^{1,2}, Necla Koçhan³, Ahmet Sezgin⁴, Gözde Ertürk Zararsız^{2,5}, Erdem Karabulut⁶, Gökmen Zararsız^{2,5}

¹Erciyes University Rectorate, Institutional Data Management and Analytics Units, Kayseri, 38280, Türkiye

² Erciyes University, Drug Application and Research Center (ERFARMA), Kayseri, 38280, Türkiye

³Izmir Biomedicine and Genome Center (IBG), Izmir, 35140, Türkiye

⁴ Abdullah Gül University, Faculty of Engineering, Department of Computer Engineering, Kayseri, 38280, Türkiye

⁵ Erciyes University, Faculty of Medicine, Department of Biostatistics, Kayseri, 38280, Türkiye

⁶ Hacettepe University, Faculty of Medicine, Department of Biostatistics, Ankara, 06140, Türkiye

e-mail: gokmenzararsiz@erciyes.edu.tr

Predicting survival in cancer patients using high-dimensional genomic or gene expression data, such as RNA-sequencing (RNA-seq), attracted much attention in recent years. Finding a correlation between time-to-event and gene expression profile of patients, for instance, can lead to more precise prognosis predictions and, consequently, better treatment strategies. Although traditional approaches have been used to analyze classical censored survival data, more sophisticated approaches, such as regularized Cox methods, or machine learning approaches adapted to survival analyses, are used to predict the patient survival times on large-scale gene expression data characterized by high dimensionality, heterogeneity, and high-collinearity (i.e., have highly correlated genes). However, due to the time and status variables in the survival data, the researcher cannot be directly applied classification algorithms with high performance and accuracy to survival analyses. In this study, we aimed to develop a new approach to analyze RNA-seq survival data by transforming data including time and status variables into new data including binary classification variables. We developed the new approach using ten RNA-seq data from various types of cancer downloaded from the TCGA database. All analyses were carried out by R language programming software. The concordance index (Harrell's c-index) and Area Under Curve (AUC) were used as a measure of performance. The results demonstrated that the newly developed approach performs as well as or better than other survival algorithms.

Keywords: Survival, RNA-seq, Machine-learning

OP15. ROBUST PROTEIN CO-EXPRESSION NETWORK FOR COVID-19

Ayca Olmez¹, Aylin Alın²

¹*The Graduate School of Natural and Applied Science, Department of Statistics,
Dokuz Eylul University, İzmir, Türkiye*

²*Faculty of Science, Department of Statistics, Dokuz Eylul University*

e-mail: olmezayca@gmail.com

COVID-19 pandemic has resulted in millions of cases and deaths worldwide, straining healthcare systems and causing significant morbidity and mortality. It has also highlighted the importance of public health measures such as vaccination campaigns, testing, contact tracing, and quarantine measures to control the spread of the virus. For making progress in these areas, understanding the defense mechanism of the virus is crucial for which Biological networks can be utilized. One of the most popular approaches in the literature is the Weighted Gene Co-expression Network Analysis (WGCNA), used to build the network structure and identify significant biological units or modules. However, WGCNA can be affected by commonly presented characteristic features in genomic datasets, such as missing values, outliers, and leverage points since it is based on Pearson correlation for similarity measure. To calculate less sensitive similarity measure to outlier observations, a robust approach based on calculating the Biweight Midcorrelation (BICOR) has previously been developed. Nevertheless, the solution to these characteristic problems in Biological network data has not been fully resolved. To deal with these issues, we propose a Partial Robust M Regression (PRM) based approach for building a robust COVID-19 protein network using protein expression values obtained from the Gene Expression Omnibus open-source database. The PRM based approach uses robust weights for potential outliers and leverage points in the data, thus improving the structure of the protein network.

Keywords: Covid-19, Protein co-expression network, Weighted gene co-expression network analysis, Partial robust m regression

OP16. infoget4gene: A USER-FRIENDLY WEB APP FOR GENETIC DATA ANALYSIS USING R SHINY

Hamdi Furkan Kepenek¹, İrem Kahveci¹, Dinçer Goksülük²

¹*Abdullah Gül University Molecular Biology and Genetics*

²*Erciyes University School of Medicine*

e-mail: hfurkan.kepenek@biogenr.com

The accessibility of open-source packages and databases, which can be utilized by researchers without advanced programming knowledge, is essential for accelerating scientific research. The R Shiny package, a popular free software in data science, provides a powerful and user-friendly tool for developing interactive web apps directly from R in various research fields, including bioinformatics, genetics, and clinical research. In this work, we present a web application built with R Shiny, which utilizes the "gget" package written in Python. Our application offers a single interface for researchers to access and analyze genetic datasets across a broad range of levels, from organism to variant analysis, without the need for coding skills. The app incorporates the core functionality of the "gget" package in Python, enabling researchers to query genomic databases with a single line of code. In addition, our application provides researchers with the ability to visualize their data and output their results in various formats, such as .png, .obj, and .pdf, depending on their desired use case. Overall, our application aims to assist researchers in analyzing their data and obtaining insights, thereby facilitating scientific research in biostatistics, bioinformatics, and data science.

Keywords: R shiny package, gget package, Web application, Genetic datasets, Data analysis

OP17. PARAMETRIC BOOTSTRAP BASED SIMULATION ON IDENTIFICATION OF DIFFERENTIALLY EXPRESSED GENES: WHICH ONE OF BORUTA OR ELASTIC NET PERFORMS BETTER?

Merve Kasikci¹, Ozgur Saman¹, Osman Dag¹

¹Department of Biostatistics, School of Medicine, Hacettepe University, Ankara, Türkiye

e-mail: mervekasikci@hacettepe.edu.tr

The identification of differentially expressed genes (DEGs) provides important insights into disease mechanisms. Machine learning algorithms are effective in discovering DEGs because they take into account relationships between genes and can process large amounts of data. Kasikci and Dag (2023) state that Determan's optimal gene selection algorithm with elastic net generalized linear models outperforms the biosigner algorithm, GMDH-type neural network algorithm, Determan's optimal gene selection algorithm with random forest, and support vector machines in terms of gene selection. In this study, we aim to compare the gene selection performances of elastic net and Boruta algorithms. Therefore, we conduct a simulation study using parametric bootstrap method. The simulation study is based on real gene expression datasets obtained from The Cancer Genome Atlas. The number of observations, number of genes, parameter estimations, and class ratios are used in the simulation scenarios. In each scenario, 10 genes are simulated as DEGs. Filtering, normalization, transformation, and univariate analysis are applied to the simulated raw datasets. Following the pre-processing step, gene selection is carried out using eight different machine learning methods. Seven of these methods include Boruta algorithm with different importance measures (Kursa and Rudnicki, 2010). The other method is Determan's optimal gene selection algorithm with elastic net generalized linear models (Determan, 2015). A confusion matrix with actual DEG and non-DEG numbers in the columns and estimated DEG and non-DEG numbers in the rows is used to assess the gene selection performances of algorithms. The results of the simulation study indicate that elastic net algorithm performs well in terms of precision. On the other hand, Boruta algorithms with random ferns and extra trees outperform the elastic net with regard to recall.

Keywords: Machine learning, Boruta algorithm, RNA-seq data, Feature selection

References:

- Kasikci, M., Dag, O. (2023). Machine Learning Based Gene Selection Algorithms for RNA-seq Data. [Manuscript submitted for publication].
- Kursa, M. B., Rudnicki, W. R. (2010). Feature selection with the Boruta package. *Journal of statistical software*, 36, 1-13.
- Determan Jr, C. E. (2015). Optimal algorithm for metabolomics classification and feature selection varies by dataset. *International journal of biology*, 7(1), 100.

OP18. DETERMINATION OF INTRON RETENTION IN GASTRIC CANCER RNA-SEQUENCE DATA BY IRFINDER-S BIOINFORMATICS ALGORITHM

Esma Gamze Aksel¹, Vahap Eldem², Selim Can Kuralay², Gökmen Zararsız³

¹Erciyes University, Faculty of Veterinary Medicine, Department of Genetic, Kayseri, Türkiye

²Istanbul University, Faculty of Science, Department of Biology, Division of Zoology, İstanbul, Türkiye

³Erciyes University, Faculty of Medicine, Department of Biostatistics, Kayseri, Türkiye

e-mail: gamzeilgar@erciyes.edu.tr

The aim of this study is the determination of genes involved in the mechanism of intron retention in gastric cancer RNA-sequence data with the IRFinder-S algorithm. For this purpose, RNA-seq data of 11 sample tumors and normal tissue biopsies located near the tumor were examined by examining open access data from the NCBI SRA database. RNA-sequencing data from control and tumor biopsies were obtained in paired-end, fastq format. Analyzes were carried out on the virtual server of the university. As the reference genome, GRCh38-Release109 of the human genome was indexed with the STAR program. After the reference genome indexing, an index was created in IRFinder-S format. Retained introns were calculated using the BAM mode in the IRFinder-S algorithm. Statistical significance of intron retention between tumor and control groups was determined with the DESeq2 package. According to the results obtained, statistical significance was found between RNA-sequence data of gastric control and tumor tissues in terms of 91 genes ($p_{adj} < 0.05$). In particular, the determination of intron retention at the transcriptome level of the cells and tissues of gastric cancer and other cancer types that can be examined is required. Recognition, use and validation of the algorithm and its types related to the determination of intron retention are important in the understanding and treatment of the cancer formation pathway. *This study was supported by Erciyes University Scientific Research Coordination Unit with project number TYL-2022-12231.*

Keywords: Gastric cancer, IRFinder-S, Intron retention, RNA-sequencing, Transcriptome

OP19. FEATURE EXTRACTION AND BIOMARKER ANALYSIS FOR DIFFERENTIATING COLON POLYPS FROM COLONOSCOPIC IMAGES

Refika Sultan Doğan¹, Ebru Aker², Serkan Doğan³, Bülent Yılmaz⁴

¹*Department of Biogineering, Abdullah Gül University, Kayseri, Türkiye*

²*Pathology Clinic, Kayseri City Hospital, Kayseri, Türkiye*

³*Gastroenterology Clinic, Kayseri City Hospital, Kayseri, 38080, Türkiye*

⁴*Department of Electrical Engineering, Gulf University for Science and Technology, Mishref, Kuwait*

e-mail: refikasultan.dogan@agu.edu.tr

Colon polyps are a common precursor to colorectal cancer. Their early detection and differentiation a crucial task in preventive medicine. Machine learning techniques are substantial for automated differentiation and characterization of colon polyps from endoscopic images and videos. Feature extraction techniques can be used to extract features that distinguish polyps from these images. This study aims to differentiate neoplastic and non-neoplastic colon polyps using feature extraction methods. The dataset consisted of 1356 video frames from 82 patients with colon polyps obtained from the Kayseri City Hospital. 51 polyps belong to neoplastic and rest of them is the non-neoplastic and, they have 844 and 512 frame number respectively. The Local Binary Pattern and Haralick features methods were used to extract 59 and 13 features, respectively, from the colon polyps. Among these features, the top 10 features that have the strongest relationship in distinguishing neoplastic and nonneoplastic groups were selected and the potential of these features as biomarkers was investigated. Chi-square statistical test was used to select the best 10 features. The biomarker performance of these best features evaluated by receiver operating characteristic analysis. The results showed that the Angular Second Moment feature had a sensitivity of 87.0% (+/-3.00) and a specificity of 74.0% (+/-4.00) in distinguishing neoplastic and non-neoplastic polyps. These results suggest that the Angular Second Moment feature could be a useful feature for distinguishing neoplastic and non-neoplastic colon polyps. Therefore, this feature has potential as a diagnostic marker for colon polyps.

Keywords: Machine learning, Feature extraction, Colonoscopy, Colon polyps, Texture analysis

OP20. WRSmoonRF: WEIGHTED ROBUST SUFFICIENT M-OUT-OF-N REGRESSION FOREST

Aylin Alın¹

¹*Department of Statistics, Dokuz Eylul University, Türkiye*

e-mail: aylin.alin@deu.edu.tr

Random Forest method is a highly data-adaptive machine learning tool. It allows for considering both regression and classification problems. In this study, we focus on the case the response variable is measured at least in interval scale, and the method is called "Random Forest Regression". It is an ensemble learning method that is a committee of decision trees from which final predictions are aggregated. The method is capable of handling large data sets, missing values in predictors, multicollinearity, or data sets where the sample size is much smaller than the number of predictors. However, it is sensitive to the outlying data points (outliers) which can be caused by human error during data collection or recording, by the participants that purposefully report incorrect data, by selecting data points from a different population than the rest of the sample or by a contaminated error distribution. Whatever the cause, outliers profoundly affect the results of an experiment. They generally increase the error variance and reduce the power of the test, altering the odds of making both Type I and II errors. All these problems can be prevented by using robust methods. Hence, we develop an approach to increase its robustness against the outliers. We implement robust weights within the process of building each tree in the forest, and a different bootstrap technique for bagging that is more robust requiring fewer data than traditional bootstrap yet giving the same and even better prediction performance. We also introduce weights for each tree in the forest. We empirically investigate the prediction performance of the proposed method on contaminated and uncontaminated simulated data sets as well as on real data. We present the theoretical background of why the proposed alternative bagging works.

Keywords: Bootstrap, Outliers, Random forest, Regression, Robustness

OP21. OPTIMIZING NUMBER OF HIDDEN LAYER AND HYPERPARAMETERS OF DEEP NEURAL NETWORK BY BAYESIAN OPTIMIZATION

Yasin Görmez¹, Duygu Korkmaz Yalçın², Sıddık Keskin²

¹*Sivas Cumhuriyet University, Faculty of Economics and Administrator Science*

²*Van Yüzyüncü Yıl University, Medicine Faculty*

e-mail: yasingormez@cumhuriyet.edu.tr

This study aims to use the Bayesian Optimization Technique proposed for methods with a large number of hyperparameters, in the optimization of the number of hidden layers and hyperparameters of the artificial neural network prediction model, which has 6 different hyperparameters. Material of this study is prostate cancer data obtained from the Kaggle dataset resource and consist of 100 observations and 10 variables (Radius, Texture, Perimeter, Area, Smoothness, Compactness, Diagnosis Result, Symmetry, Fractal Dimension and Outcome). The dataset was randomly divided into 70% training and 30% test set to be used to train and test the final model. Then, the training data set was randomly divided into 80% training and 20% test set to be used for optimization. After the datasets were generated, an artificial neural network model was developed using the Python Keras library and the parameters of the model were optimized using the training dataset for testing and optimization. Type of optimized hyperparameters (Number of Dense Layer, Learning Rate, Dropout Rate, Number of Hidden Neuron in Dense Layer, Epochs, Batch) and types of values determined for these hyperparameters (Integer, Real, Categorical), ranges assigned for value spaces, and optimum values were determined. After determining the optimum number of hidden layers and hyperparameters, two different models were trained and tested using training and test data. In the first of these models, the default hyperparameters and in the second, the optimized hyperparameters were used and the performance criteria of the models were compared. Parameters in the optimized model as to default model as follows respectively. Accuracy increased from 76.66% to 83.33%, Precision from 70.00% to 77.77%, F1-Score from 80.00% to 84.84%, MCC from 0.57 to 0.68. And Recall was calculated as 93.33% in both models. Consequently, model performance has increased in the optimized model with Bayesian optimization compared to the default model.

Keywords: Bayesian optimization, Hyperparameter, Machine learning, Neural networks

OP22. AUTOMATED MACHINE LEARNING APPROACH IN CLINICAL SETTINGS: PREDICTING THE FUTURE RISKS

Didem Turgut^{1,2}, Deniz Ilhan Topcu³, Samet Senel⁴, Cuneyt Ozden⁴

¹Ankara City Hospital, Division of Nephrology, MD

²Hacettepe University, Department of Biostatistics, MSc

³Izmir Tepecik Education and Research Hospital, Department of Biochemistry, MD

⁴Ankara City Hospital, Department of Urology, MD

e-mail: dr.didem@gmail.com

Kidney stone disease is a common problem in primary care practice. Complicated stone formation can also lead to progressive loss of kidney functions at a young age. The aim of the present study was to investigate the risk factors of kidney stones in adult population and estimate the disease progression risk. While estimating these risks we used H2O automated machine learning (AutoML) tool. Boruta algorithm was used for feature selection. We conducted a retrospective cohort study. Age, sex, number of stones, stone location, surgery type, presence of hydronephrosis preoperatively, existing comorbidities, residual stone status after operation, pre-, and postop leucocyte count, and preop eGFR measurements were selected for postop eGFR prediction at the first month with machine learning (ML). 300 patient records were filtered that contain all parameters. Data set was split into training (n = 200) and test sets (n =100). H2O AutoML engine was used to develop multiple ML models for eGFR prediction. The best model was determined using the R² performance metric and then model performance was reevaluated using the test set. A total of 10 ML models were developed by the H2O engine developed. The generalized linear model had the highest R² value for the training set (0.74) and R²= 0.70 for the test set. The most three important variables for prediction of postoperative first month eGFR were, preop eGFR (31%), age (10%), and patients with multiple comorbidities (%4). Kidney stones are a risk factor for chronic kidney disease. ML algorithms are novel strong prediction tools for these progressive clinical problems.

Keywords: Machine learning, Kidney stone, Prediction

OP23. BIOINFORMATICS AND BIOSTATISTICAL MODELS FOR ANALYSIS AND PROGNOSIS OF ANTIMICROBIAL RESISTANCE

Maya Zhelyazkova¹, Stefan Tsonev², Dimitar Vassilev³

¹*Sofia University, Faculty of Mathematics and Informatics, Sofia, Bulgaria*

²*AgroBioInstitute, Sofia, Bulgaria*

³*Sofia University, Faculty of Mathematics and Informatics, Sofia, Bulgaria*

e-mail: zhelyazkova@fmi.uni-sofia.bg

The bioinformatics and biostatistics models revealing the relationship between bacteriophages and antimicrobial resistance could be regarded as important part in the scope of studying the potential impact of microbiology on contemporary medicine and pharmaceuticals. Till now research projects have some ambiguous outcomes by examining and confirming the influence of bacteriophages on variation of antimicrobial resistance genes. The major goal of our study is to apply new bioinformatics and biostatistical models in order to acquire new knowledge how the viruses, their hosts and antimicrobial resistance genes are related and how these relations can be clarified in the context of antimicrobial resistance dissemination through phages. The presented work is oriented towards zooming in the relationship and possible dependencies between bacteriophages and antimicrobial resistance. The data has been collected from different city environments all over the world. Our analyses consist of different bioinformatics and biostatistical methods for assessment of the differential abundance of phages, their diversity across samples, the impact on antimicrobial resistance categories and associations with antimicrobial resistant genes.

Keywords: Antimicrobial resistance, Bacteriophages, Diversity indexes, Relative risk, Bayesian spatial analysis

OP24. DEEP NEURAL NETWORKS FOR AVERAGE TREATMENT EFFECT ON BIOLOGICAL NETWORKS

Mehmet Ali Kaygusuz¹, Vilda Purutçuoğlu

¹*Anadolu University, Department of Economics, 26780, Türkiye*

²*Middle East Technical University, Department of Statistics, 06800, Türkiye*

e-mail: makaygusuz1988@gmail.com

The estimation of heterogeneous treatment effect is a crucial problem in machine learning, bioinformatics and deep learning, and the average treatment effect (ATE) basically measures the difference in mean (average) outcomes between units assigned to the treatment as well as units assigned to the control in such a way that the bias can be defined. Moreover, causal forests, which are an adaptation of the random forest algorithm, have a lot attention in recent years in the problem of heterogeneous treatment effect estimation. This task can be computationally very tractable, in particular, for clustered data when we work on high number of parameter (p) regarding number of observations (n). On the other hand, more recently, feedforward neural networks have been applied in different fields, by binding the random forest with deep neural network (DNN) in such a way that the random forest can detect the features in DNN and DNN can use these features to improve the performance of classification when $n \ll p$. Hereby, in this study, we aim to add the detection of the average treatment effect via random forest inserted DNN. Then, we compare the findings under with and without DNN. By this way, we consider to investigate the gain in accuracy by including DNN in average treatment effect estimation via random forest. We examine proposed algorithm with the different n and p to assess its efficiency.

Keywords: Heterogeneous effect, Random forest algorithm, Artificial neural networks, Biological data

OP25. SIMULTANEOUS SCORING OF CLUSTERS IN RECURSIVE CLUSTER ELIMINATION, APPLIED ON TRANSCRIPTOMIC DATA ANALYSIS

Nurten Bulut¹, Burcu Bakir-Gungor², Bahjat F. Qaqish³, Malik Yousef⁴

¹*Department of Computer Engineering, Abdullah Gul University, Kayseri 38080, Türkiye*

²*Department of Computer Engineering, Abdullah Gul University, Kayseri 38080, Türkiye*

³*Department of Biostatistics, University of North Carolina at Chapel Hill, NC, Chape Hill, USA*

⁴*Department of Information Systems, Zefat Academic College, Zefat 13206, Israel*

e-mail: nurten.bulut@agu.edu.tr

Gene expression data analysis is a challenging task due to the small sample size and high dimensionality of data. Feature selection (FS) is a useful approach to deal with high dimensionality. Support Vector Machines - Recursive Cluster Elimination (SVM-RCE) is one of the few techniques created for FS in high-dimensional data. SVM-RCE approach has been utilized for detecting distinct clusters of genes that exhibit varying levels of expression in the presence of pathological states. In the original study, the genes were grouped into clusters utilizing the K-means algorithm. For each cluster, a corresponding sub-data is extracted, where this sub_data contains gene expression values just for those selected genes, retaining the original class labels of the samples. Then, a score was assigned to each of these clusters by running SVM and by applying internal cross-validation. Following that, a percentage (e.g., 10%) of clusters with low scores are eliminated. The elimination is performed by filtering out the genes that are members of those lower-scored clusters. The aforementioned procedure is iteratively executed until a specific number of clusters persists. The present study proposes a novel approach for scoring the clusters simultaneously by applying linear SVM or Random Forest (RF) on the clusters centers. The absolute values of the coefficients given by linear SVM serves as the score of the corresponding cluster. Similarly, the feature weights given by RF serve as the scores of each corresponding cluster. In our experiments, we used 17 transcriptomic datasets, obtained from GEO. The results indicate that the performance of the new approach is comparable to that of the original approach. The novel methodology achieves a task completion rate that is 80% faster than the original approach. The techniques devised will facilitate disease diagnosis, in addition to enhancing our comprehension of the molecular mechanisms underlying disease onset and progression.

Keywords: Feature selection, Clustering, Recursive cluster elimination, Gene expression data analysis, Transcriptomic

OP26. THE G-S-M, GROUPING, SCORING AND MODELING APPROACH. APPLICATION OF BIOLOGICAL DOMAIN KNOWLEDGE FOR GROUPS SELECTION ON GENE EXPRESSION DATA

Malik Yousef¹, Burcu Bakir-Gungor²

¹*Department of Information Systems, Zefat Academic College, Zefat 13206, Israel*

²*Department of Computer Engineering, Faculty of Engineering, Abdullah Gul University, Kayseri 38080, Türkiye*

e-mail: malik.yousef@gmail.com

In the last two decades, there have been massive advancements in the high throughput technologies, which resulted in exponential growth of public repositories of gene expression datasets for various phenotypes. It is possible to unravel biomarkers by comparing the gene expression profiles in different conditions (e.g. cases vs. controls, drug treatments after different time points, different tissues). This problem refers to a well-studied problem in the machine learning domain, i.e., the feature selection problem. In biological data analysis, most of the computational feature selection methodologies were taken from other fields, without considering the nature of the biological data. Thus, integrative approaches are necessary for this kind of data. The main aim of integrative gene selection is to generate a ranked list considering both statistical metrics applied on the gene expression data, and the biological background information provided in external databases or through omics data sets. During the last years, we have developed an integrative approach that performs groups selections rather than feature selection. The generic approach is called G-S-M (Grouping, Scoring, and Modeling). I will present our related works in this topic such as SVM-RCE, maTE, CogNet, miRcorrNet, miRModuleNet, 3Mint, GediNET, PriPath, TextNetTopics, GeNetOntology, MicroBiomeNet, mirDisNet, and other tools under development.

Keywords: Integrates prior biological knowledge, Gene expression, Machine learning, Feature selection

OP27. THE EFFECT OF MISSING DATA IMPUTATION METHODS ON CLASSIFICATION PERFORMANCE ACCORDING TO DIFFERENT MISSING RATES IN HIGH DIMENSIONAL DATA

Buğra Varol¹, İmran Kurt Ömürlü², Mevlüt Türe²

¹*Adnan Menderes University, Institute of Health Sciences, Division of Biostatistics, Aydın*

²*Adnan Menderes University, Faculty of Medicine, Division of Biostatistics, Aydın*

e-mail: bugravarol87@gmail.com

Missing data is an important problem in the analysis and classification of high-dimensional data. The aim of this study is to compare the effects of six different missing data imputation methods on classification performance in high-dimensional data. In this study, missing data imputation methods were evaluated using data sets, whose independent variables are mixed correlated with each other, for binary dependent variable, $p=500$ independent variables, $n=150$ units and 1000 times running simulation. Missing data structures were created according to the missing at random (MAR) mechanism and different missing rates. Different datasets were obtained after having imputed the missing values separately by six imputation methods including ridge regression, lasso regression, elastic net regression, support vector machine (SVM), extreme gradient boosting (XGBoost), and extreme learning machine (ELM). At the end of the simulation, classification scores of the methods by gradient boosting machines (GBM) were obtained and the missing value prediction performances were evaluated according to the distance of these scores from the full data sets. The increase in the missing rate affects the classification performance in high-dimensional data.

Keywords: Missing data, Simulation, Imputation, Classification, Gradient boosting machines

OP28. HEALTH SPACE MODEL USING DEEP LEARNING

Taesung Park¹, Chanhee Lee²

¹*Department of Statistics, Seoul National University, Seoul 08826, Korea*

²*Interdisciplinary Program in Bioinformatics, Seoul National University, Seoul 08826, Korea*

e-mail: tspark@stats.snu.ac.kr

As the era of personalized medicine arrives, measuring and visualizing an individual's health status in an objective way is becoming increasingly crucial. To visualize high-dimensional data, many dimension reduction techniques have been developed in the fields of statistics and machine learning, such as Principal Component Analysis (PCA), t-distributed stochastic neighbor embedding (t-SNE), and uniform manifold approximation and projection for dimension reduction (UMAP). However, the axes created from these methods do not have a clear biological interpretation. To improve biological interpretability, a more informative statistical method called 'health space (HS)' has been developed using statistical models such as the logistic regression model and the proportional odds model. However, this HS model lacks flexibility in modeling non-linear relationships between the phenotype and independent variables. We developed a new HS model named Deep Ordinal Neural Network (DONN) to model non-linear relationship and utilize ordinal information of the phenotype. DONN achieves these goals by combining Deep Neural Network (DNN) structure and cumulative logit with shared coefficient values. We trained DONN using 32,140 samples from Korea National Health and Nutrition Examination Survey (KNHANES) and then validated the model using Ewha-Boramae cohort data with 862 samples and the Korea association resource project data with 3,199 samples. Health status of an individual was collected as ordinal phenotype (healthy, risk 1, risk 2, disease). The proposed DONN model was compared with existing health space models based on proportional odds model and binary DNN models. Both training and external datasets showed that DONN had best performance in discriminating health status. DONN models non-linear relationship and utilizes ordinal information of the phenotype well. DONN can be an effective tool for visualizing an individual's health status in an objective way.

Keywords: Visualization, Health space, Deep ordinal neural network

OP29. HOW TO EXPLAIN CARBOHYDRATE METABOLISM DISORDERS USING MACHINE LEARNING MODELS?

Deniz İlhan Topcu¹, Banu Isbilen Basok¹

¹*University of Health Sciences Izmir Tepecik Training and Research Hospital
Department of Medical Biochemistry*

e-mail: ditopcu@gmail.com

It is known that using machine learning (ML) models to improve patient care is a novel approach, but its implementation in clinical practice remains limited due to the black-box nature of ML models. In medicine, making clinical decisions requires a more explainable approach, such as "white box" statistical processes. Various explainability tools promise a better understanding of developing ML models and, hence, the implications of those. In this study, we aim to evaluate the best-fitting ML models regarding not only their performance metrics but also novel explainability results. To evaluate different aspects of carbohydrate metabolism, two separate datasets were used in the study: (1) 3,036 records from the NHANES open dataset for hemoglobin A1c (A1c) classification (healthy, prediabetes, and diabetes) and (2) 226 actual patient data from a tertiary hospital to classify impaired glucose tolerance (IGT). For A1c classification, two tree-based ML algorithms were developed, and global and local model-agnostic explanation methods, including performance metrics, feature importance, partial dependence, and Shapley additive explanation plots (SHAP), were carried out. To classify IGT, Random Forest and LightGBM models were created, and SHAP and permutation importance analyses were performed to determine feature importance. Our results showed that, for A1c classification, while both models had similar performance metrics, they exhibited slightly different variable importance and local explainability results. The same trend was also observed for classifying IGT. The global explainability results within and between the two ML models were quite different. In this study, explainable AI models for two different clinical aspects of carbohydrate metabolism were evaluated using open-source and real patient datasets. The findings suggest that, despite existing constraints on explainability approaches, especially for clinical interpretation, a combination of global and local explanation models provides insights into model evaluation and can be employed for improving or contrasting different models.

Keywords: Artificial intelligence, Carbohydrate metabolism, Explainability, ExAI, Machine learning

OP30. ABNORMALITY DETECTION AND CLASSIFICATION ON MAMMOGRAPHY IMAGES VIA CONVOLUTIONAL NEURAL NETWORKS

Hanife Avci¹, Gamze Durhan², Figen Demirkazık², Meltem Gülsün Akpınar², Jale Karakaya¹

¹*Department of Biostatistics, Hacettepe University*

²*Department of Radiology, Hacettepe University*

e-mail: hanife.avci@hacettepe.edu.tr

Breast cancer is the most common type of cancer among women. Digital mammography screening is an imaging method that helps determine the rate of early detection of breast cancer. In the last two decades, Computer Aided Detection (CAD) systems were developed to help radiologists analyze screening mammograms. Since 2012, deep convolutional neural networks (CNN) have been shown to achieve high performances in image recognition. In this study, the effectiveness of deep learning using Convolutional Neural Networks (CNN) was tested on digital mammography images in detecting the presence of lesions (abnormal/normal). In addition, classification performances of different CNN architectures will be compared. Digital mammography images were used for this study, which received ethics committee permission from Hacettepe University. Craniocaudal (CC) view were carried out for each breast. Spyder environment (with the Keras and TensorFlow) was used for image processing steps (pre-processing, segmentation, feature selection, classification) and evaluating the performance on the classification. The data set was divided into two as training and test sets with 80% and 20%. Noise was removed from the images with image pre-processing techniques. Then, the regions of interest (ROI) obtained after segmentation were determined by CNN with linear filters and activation functions. The features were obtained from the ROIs with the help of the GLCM matrix. So far, the performance of 87 lesioned image (malign and benign images) 59 lesion less images (normal image) in the CNN model was examined. The applied CNN model reached 79.64%, 80.85% and 75.44% performance rates for accuracy, sensitivity and specificity, respectively. By increasing the number of observations, the analyzes will be renewed and also the classification performances of different CNN architectures will be compared.

Keywords: Mammography classification, Deep learning, Medical imaging processing, Computer-aided detection

OP31. ESTIMATING THE COST OF COVID-19 TO TURKISH TOURISM WITH TIME SERIES AND MACHINE LEARNING MODELS

Günel Bilek¹

¹*Department of Business Administration, Izmir Democracy University, Izmir, Türkiye*

e-mail: gunalbilek@gmail.com

This study aims to predict the cost of COVID-19 to the Turkish tourism sector with time series and machine learning approaches. The study data comprises monthly number of tourists visiting Türkiye between 2003 and 2022. The fitted models are SARIMA, NNAR, TBATS, ETS and SARIMA-NNAR, SARIMA-TBATS, SARIMA-ETS hybrid models. The model performances were compared by error metrics mean error (ME), mean absolute error (MAE), root mean square error (RMSE), and mean absolute percentage error (MAPE) calculated on new data and the model with the minimum error metrics were chosen as the best model. SARIMA(1, 1, 0)(1, 1, 0)₁₂ performed best among all the candidate models and passed model diagnostics checks, making it the best model. Furthermore, the chosen model was used to forecast the number of visitors for the years 2020, 2021 and 2022 and these numbers were compared to the observed values of the same period in order to detect the effect of COVID-19 on the Turkish tourism sector. According to our model, 103.482.894 less people visited Türkiye due to COVID-19 in 2020, 2021 and 2022, corresponding to a loss of 92.200.021.042 USD to Turkish tourism. Finally, unlike the SARIMA (1, 1, 0)(1, 1, 0)₁₂ model, the performances of hybrid and machine learning models were not satisfying. They succeeded in estimating the seasonality, but failed in estimating the trend.

Keywords: Time series, SARIMA, NNAR, TBATS, ETS, Hybrid models

OP32. ASSESSING PREDICTION ACCURACY OF JOINT MODELS: A NOVEL APPROACH BASED ON MUTUAL INFORMATION CRITERION

Merve Basol Goksuluk¹, Dincer Goksuluk¹, A. Ergun Karaagaoglu²

¹*Department of Biostatistics, Faculty of Medicine, Erciyes University, 38280 Kayseri, Türkiye*

²*Department of Biostatistics, Faculty of Medicine, Lokman Hekim University, 06510, Ankara, Türkiye*

e-mail: dincergoksuluk@erciyes.edu.tr

Joint modeling has emerged as a preferred approach in follow-up studies to analyze the relationship between longitudinal and time-to-event data. This approach allows for dynamic and personalized predictions of disease progression for individual patients. However, the accuracy of these predictions needs to be evaluated to ensure their reliability. Common methods for model evaluation include the area under the receiver operating characteristic (ROC) curve and the Brier score. In this study, we propose a novel approach to interpret the performance of joint models using mutual information criteria. Our proposed approach provides insights into the amount of information gained from each time point or the optimal timing for measurements after the last one, allowing us to answer questions such as which time point provides the most information or when to take measurements after the last one. However, the latter aspect is not covered in this proposal. We applied our approach to a real-world dataset on peritoneal dialysis, specifically examining the relationship between serum albumin levels and mortality using joint models. We then compared the prediction accuracies of our approach with other methods through simulation studies under various scenarios. Our findings suggest that mutual information criteria can be used alongside traditional methods for model evaluation. Simulation study results indicate that increasing sample size improves model performance, but the number of repeated measurements after a certain time point does not significantly affect the results. In conclusion, mutual information criteria can be a valuable tool to assess the prediction accuracy of joint models in conjunction with other methods. While our approach is not claimed to be the best, it can complement existing approaches for model evaluation.

Keywords: Longitudinal, Time-to-event data, Joint model, Accuracy, Mutual information

OP33. PRACTICAL CONSIDERATIONS FOR STATISTICAL MODELS AND THEIR IMPLEMENTATIONS OF PHASE I DOSE-ESCALATION ONCOLOGY TRIALS

Burak Kürsad Günhan¹, Pavel Mozgunov², Anja Victor¹

¹*Merck Healthcare KGaA, Darmstadt, Germany*

²*MRC Biostatistics Unit, Cambridge, UK*

e-mail: burak-kuersad.guenhan@merckgroup.com

Phase I dose-escalation trials in oncology constitute the first step in investigating the safety of potentially promising therapies in humans. Conventional statistical methods for these trials include algorithm-based designs, such as 3+3 rule, which have been developed for chemotherapy treatments. In recent years, there is an increasing need for new statistical developments, since trials have become more complex with the move to new types of treatments like targeted agents and immunotherapy as well as the increased interest in combination of different therapies. This talk will give an overview over new statistical developments in oncology dose escalation with focus on testing different treatment schedules and combination. Starting from introducing the application of Bayesian logistic regression model (Neuenschwander et al., 2008), this talk will move on to present how to handle a change in regimen by adding a covariate into the model (Bailey et al., 2009). Finally, the extension for changes in regimen in a combination setting will be showcasing the implementation of the extended Partial Ordered Continual Reassessment Method (Mozgunov et al., 2022). Special focus will be on presenting the cross-institution collaboration for the development of the R package **crmPack** (Sabanés Bové et al., 2019 and Sabanés Bové et al., 2022), that is working on implementing these new methods in a user friendly way.

Keywords: Bayesian methods, Dose-escalation, Phase I oncology trial, R

References

- Neuenschwander, B., Branson, M., and Gsponer, T. (2008). Critical aspects of the Bayesian approach to phase I cancer trials. *Statistics in Medicine*, 27(13), pp.2420-2439. doi: [10.1002/sim.3230](https://doi.org/10.1002/sim.3230).
- Bailey, S., Neuenschwander, B., Laird, G., Branson, M. (2009). A Bayesian case study in oncology Phase I combination dose-finding using logistic regression with covariates. *Journal of Biopharmaceutical Statistics*, 19(3), 469-84. doi: [10.1080/10543400902802409](https://doi.org/10.1080/10543400902802409).
- Mozgunov, P., Jaki, T., Gounaris, I., Goddemeier, T., Victor, A., and Grinberg, M. (2022). Practical implementation of the partial ordering continual reassessment method in a phase I combination-schedule dose-finding trial. *Statistics in Medicine*, 41(30), 5789- 5809. doi:[10.1002/sim.9594](https://doi.org/10.1002/sim.9594).
- Sabanés Bové, D., Yeung, W. Y., Palermo, G., and Jaki, T. (2019). Model-Based Dose Escalation Designs in R with crmPack. *Journal of Statistical Software*, 89 (10), 1-22. doi: [10.18637/jss.v089.i10](https://doi.org/10.18637/jss.v089.i10).
- Sabanés Bové, D., et al. (2022). Improving Software Engineering in Biostatistics: Challenges and Opportunities. Under review, Preprint: doi: [10.48550/arXiv.2301.11791](https://doi.org/10.48550/arXiv.2301.11791).

OP34. THE CHERNOFF FACES METHOD FOR VISUALIZING COMPLEX DATA: AN APPLICATION FOR IDENTIFYING DIFFERENCES BETWEEN COVID-19 AND CONTROL GROUPS

Elif Kaymaz¹, Ferhan Elmalı¹, Büşra Emir¹, Fatma Ezgi Can¹, Mustafa Ağâh Tekindal¹

¹*Izmir Katip Celebi University, Faculty of Medicine, Basics Medical Science, Department of Biostatistics*

e-mail: kaymaz.elif@yahoo.com

Data visualization is a crucial topic in today's world due to its ability to present statistical and variable data in large datasets in an easily understandable manner through graphical interfaces, making it easier to comprehend the data. Chernoff faces is one of the methods that can visually express complex and multivariate data. This method creates faces by combining different organ sizes and colors that represent different dimensions of the data using caricature drawings of human faces. The sizes and positions of the organs used in drawing the faces are adjusted to reflect different features of the data. Currently, research on post-Covid-19 infection continues to exist. This epidemic has led to research in many areas, including the diagnosis, treatment, and transmission of the disease. In this study, Chernoff faces were used to determine the differences between the Covid-19 and control groups. Laboratory parameters of blood samples taken from the Covid-19 and control groups were measured. Chernoff faces were used to make these data more understandable. The results showed that the differences between the Covid-19 and control groups could be determined by using Chernoff faces. Differences in organ sizes that represent critical factors, such as blood parameters, played a significant role in determining the differentiation between the groups. Studies have shown that the human brain processes visual information 60,000 times faster than text, which is why data visualization techniques are frequently used in analysis and reporting studies, especially in the healthcare field. In conclusion, the method of using Chernoff faces to determine the differences between Covid-19 and control groups provides a more easily understandable and readable presentation compared to other data visualization methods. This study proposes a new method that can be used in analysis and reporting studies in the healthcare field.

Keywords: Chernoff Faces, Data Visualization, Statistics

OP35. A REVIEW ON RANDOMIZED CONTROLLED TRIALS IN EMERGENCY MEDICINE: METHODOLOGICAL ISSUES

Demet Arı¹, Pınar Günel², Buket İpek Berk³, İhsan Berk², Vildan Sümbüloğlu²

¹Gaziantep SANKO University, Faculty of Medicine, Department of Emergency Medicine

²Gaziantep SANKO University, Faculty of Medicine, Department of Biostatistics

³Gaziantep SANKO University, Graduate Education Institute, Department of Biostatistics

e-mail: buketipek5@gmail.com

The subject of methodological and statistical review of clinical research has been included in the literature since many years. The most important reason giving preference to this subject is to draw attention to frequently recurring errors/omissions and prevent the recurrence of them. For the validity and accuracy of the results of a clinical trial, the research method should be designed appropriately and research should be conducted exactly according to this method. Otherwise, the reliance to the results obtained from clinical trials decreases. In order to evaluate the quality of a research, all the necessary information should be included in detail in the method section. This study includes a methodological review of randomized controlled trials conducted in emergency medicine departments. For this purpose, the Pubmed database was scanned with the keywords “emergency department, pain management” and randomized controlled clinical trials published in the last five years were determined. We aimed to examine the method sections of all 125 articles listed according to these criteria. However, at the first stage, 50 randomly selected articles were considered as preliminary and the topics examined in accordance with the CONSORT check list were determined. These topics are; whether there is a biostatistics specialist in the research team, the location, date and design of the research, interventions, population, sample, control group, randomization, blinding, bias, data collection tools and plan, outcome variables, inclusion and exclusion criteria, ethics committee approval, informed written consent and statistical method. In the 50 articles examined, randomization was the topic that was found to be the most common error. Although all of the examined articles included a description of randomization, the method used in 94% of them was not appropriate. This is followed by statistical method; there is insufficient / incomplete explanation in 36% of the articles, incorrect use of method in 12%, both insufficient / incomplete explanation and incorrect use in 14%. In addition, it is noteworthy that in about all of the articles (92%) there is no biostatistics specialist or its existence is unclear. The method of blinding, which is an important issue in terms of bias, was not applied even though it was required in 10% of the articles. There were less deficiencies in the other titles examined relatively. In order to prevent methodological issues in a research, it is of great importance for journal referees to conduct method reviews more carefully and journal editors to be more sensitive on this matter.

Keywords: Emergency medicine, Randomized controlled trials, Methodological issue

OP36. NONPARAMETRIC ESTIMATION OF DISTRIBUTION FUNCTION USING RANKED SET SAMPLING WITH UNEQUAL PROBABILITIES

Yusuf Can Sevil¹, Tuğba Özkal Yıldız²

¹*The Graduate School of Natural and Applied Sciences, Dokuz Eylül University, Izmir, Türkiye*

²*Department of Statistics, Faculty of Science, Dokuz Eylül University, Izmir, Türkiye*

e-mail: yusufcansevil92@gmail.com

Ranked set sampling (RSS) is a widely used sampling technique when collecting measurements is difficult and/or costly. Estimating distribution functions from experimental data is a common problem, and several authors have proposed empirical distribution functions (EDFs) based on RSS and its modifieds. There are a few studies on distribution function estimation in finite population setting. Also, these studies have been developed classical EDF, which assigns equal weight to each sampling unit. The present work focuses on developing design-based estimators for distribution function using RSS methods (level-0, level-1, and level-2). These sampling methods have different procedures in terms of their replacement policies. Thus, different inclusion probabilities depending on the sampling methods are assigned to population units for selection. The design-based estimators use the inclusion probabilities unlike the classical EDF. Therefore, this work brings a new perspective on the estimation of the distribution function by using RSS techniques. The bias and efficiency of the proposed estimators are investigated theoretically and numerically. According to the results, level-2 sampling method shows outperformance among the RSS methods and simple random sampling. Also, we investigate a pointwise estimate of the distribution function and confidence interval of the median of sheep weight at 7 months using level-2 sampling procedure in real data application.

Keywords: Ranked set sampling; Probability sampling designs; Design-based estimators; Empirical distribution function

OP37. A COMPREHENSIVE COMPARISON OF LOW DENSITY LIPOPROTEIN CHOLESTEROL EQUATIONS

Serra İlayda Yerlitaş^{1,2}, Gözde Ertürk Zararsız^{1,2}, Halef Okan Doğan³, Serkan Bolat³,
Necla Kochan⁴, Ahu Cephe⁵, Gökmen Zararsız^{1,2}, Arrigo F.G. Cicero^{6,7}

¹Department of Biostatistics, Erciyes University, Kayseri, Türkiye
Drug Application and Research Center (ERFARMA), Erciyes University, Kayseri, Türkiye

³Department of Biochemistry, Cumhuriyet University, Sivas, Türkiye

⁴Izmir Biomedicine and Genome Center (IBG) 35340, İzmir, Türkiye

⁵Institutional Data Management and Analytics Unit, Erciyes University Rectorate 38280, Kayseri, Türkiye

⁶Medical and Surgical Sciences Dept., Alma Mater Studiorum University of Bologna 40126, Bologna, Italy

⁷IRCCS AOU S. Orsola-Malpighi di Bologna, Bologna, Italy

e-mail: serrayerlitas@erciyes.edu.tr

In recent years, cardiovascular disease (CVD) has increased worldwide. CVDs accounted for 32% of global deaths in 2019. Low-density lipoprotein cholesterol (LDL-C) associated with these diseases is an important marker used to investigate appropriate treatment strategies and is also used as a target measure in clinical practice guidelines. The gold standard for LDL-C measurement is ultracentrifugation and β -quantitation. However, direct measurement of LDL-C is too expensive for most laboratories and performs poorly at high triglyceride concentrations. Researchers have developed formulas for estimating LDL-C in the literature, but there is no consensus on the best method for estimating LDL-C. In this study, a comprehensive comparison was made of the performance of existing equations to estimate directly measured LDL-C levels in 3 different populations (Korean, Italian, and Turkish) and with 3 different homogeneous assays (Roche, Beckman, and Siemens). In addition, in line with the reference studies where lipid measurement standards may differ in the pediatric population, the equations were compared in the pediatric group on 3 different direct assays (Roche, Beckman, Siemens). The study included 278,695 adults and 3,908 children participants. 22 equations were used in the comparisons, including formulae such as Friedewald, Martin/Hopkins, and Sampson. When the results of the study were examined, it was found that the performance of the equations changed significantly as the populations and direct assays changed. Therefore, there is a need to develop a new high performance LDL-C estimation equation that has been validated in a larger population and using different assays.

Keywords: Low-density lipoprotein cholesterol (LDL-C), Lipid metabolism, Cardiovascular risk

OP38. PUBLIC HEALTH-FOCUSED USE OF COVID-19 RAPID ANTIGEN PCR TESTS

Yonatan Woodbridge^{1,2}, Yair Goldberg³, Sharon Amit⁴, Naama M. Kopelman²,
Micha Mandel⁵ and Amit Huppert^{1,6}

¹*The Gertner Institute for Epidemiology & Health Policy Research, Sheba Medical Center, Ramat Gan, Israel*

²*Department of Computer Science, Holon Institute of Technology, Holon, Israel*

³*The Faculty of Data and Decision Sciences, Technion–Israel Institute of Technology, Haifa, Israel*

⁴*Clinical Microbiology, Sheba Medical Center, Ramat Gan, Israel*

⁵*Statistics and Data Science, The Hebrew University of Jerusalem, Jerusalem, Israel*

⁶*Faculty of Medicine, Tel Aviv University, Israel*

e-mail: amit.huppert@gmail.com

During the Covid-19 pandemic, accurate PCR tests were augmented by the cheap, rapid, and logistically convenient, yet less sensitive antigen tests. This testing policy shift was implemented due to limited availability of PCR tests during the Omicron surge, but took place without proper quantification of the tradeoffs. Yet, evidence-based surveillance requires a robust understanding of the strengths and limitations of the available detection methods. Using 41,065 paired tests from this period, we estimate how the sensitivity of antigen tests changes as a function of Ct value and other key covariates. The results reveal a logarithmic relationship between antigen detection probability and viral load, as quantified by Ct-values of the PCR tests. Further analysis shows a statistically significant association with an odds ratio of approximately 0.76 with each unit of Ct-value. The analysis suggests that in spite of their compromised sensitivity, antigen tests are a natural solution for routine use, while PCR tests should be considered in situations where a false negative result could have serious consequences. These findings are the foundations of policies that will utilize the strengths of the different tests, and achieve enhanced hybrid surveillance.

Keywords: Informed decisions, PCR, Antigen, COVID-19, Logistic regression

OP39. DIABETIC RETINOPATHY DIAGNOSIS AND CLASSIFICATION

Berk Pişkin¹, Aylin Alın², Rim Khazhin³, Ahmet Mert Saygu⁴, Ahmet Ömer Özgür⁵

¹*Department of Statistics, Dokuz Eylül University the Graduate Schools of Natural and Applied Sciences*

²*Department of Statistics, Dokuz Eylül University The Faculty of Science*

^{3, 4, 5}*Eye Checkup, Antalya*

e-mail: berkpiskin.a@gmail.com

Diabetic retinopathy (DR) is an eye disease that occurs in diabetic patients. If not diagnosed early, DR can be detected at a much later stage and the patient may already have irreversible eye damage. Therefore, early diagnosis of DR and the identifying important factors of DR are key to correctly planning the treatment of DR. In this study, we utilize the data generated by feature extraction methods from the fundus image database of EyeCheckup for classification and variable importance detection, where the data is used in accordance with confidentiality policies. Support Vector Machines, Random Forest, XGBoost and Logistic Regression classifiers are selected to be employed on the data, since they have proven themselves in the literature. SMOTE (Synthetic Minority Over-Sampling Technique) and SMOTE variants will be used to handle class imbalance. Finally, permutation importance method will be used on the best performing model after the performance comparison within the proposed metric (SS-Score). In the study, our aim is to perform these analyzes on 3 variations of the data specific to each class and to determine the factors affecting each level of the disease separately. This study will contribute to DR classification studies and the field of machine learning by gradually combining methods that have not been combined in the same study before, by comparing the models separately on the basis of each metric to be examined after the models are created, by analysis and interpretation of the results, and by the positive contributions that the findings will make to the diagnosis of DR. Another contribution is that this study will serve as a resource for future researchers to solve the problems they face in classification studies.

Keywords: Diabetic retinopathy, Machine learning, Class imbalance, Feature importance, SMOTE

OP40. A NEW ESTIMATOR FOR THE DISCRIMINATION ACCURACY IN A FOUR-CLASS CLASSIFICATION PROBLEM

Elena Nardi¹

¹*IRCCS Azienda Ospedaliero-Universitaria di Bologna;
Department of Surgical and Medical Sciences, University of Bologna, Bologna*

e-mail: elena.nardi2@unibo.it

The present work focuses on the study and extension of ROC analysis methodology for multiple-class classification problems. In clinical medical research, the need for developing an approach to measure the diagnostic accuracy of a biomedical test in classifying the true status of a patient is a critical point when doing both diagnosis and prognosis. In a two-category classification setting, the ROC curve analysis is the most natural approach and the Area Under the Receiver Operating Characteristic curve (AUC) is a summary measure of the diagnostic accuracy. However, many real classification problems rely to more than two classes; consequently, the notion of the Area Under the Curve has been extended to the Volume Under the Surface (VUS) and, in the more complex situations of more than three classes, to the hypervolume (HUM). In this contribution, we develop a new estimator of the accuracy measure of a continuous diagnostic marker in a four-class classification framework. Our proposal is based on the Lehmann family ROC surfaces as in Nze Ossima et al., 2015. In particular, we derive the analytical form of the HUM estimator and the analytical representation of its variance. To assess the performance of the proposed estimator and compare it with the two alternatives existing in the literature, simulation exercises and empirical applications are presented.

Keywords: Classification, ROC analysis, Proportional hazards

OP41. REVIEW: ACCESS CONTROL IN ELECTRONIC MEDICAL RECORDS (EMR)

Hilah Alnafisah¹, Rawaby Alsaaid¹, Fatimah M. Alturkistani¹

¹*College of Computer and Information Sciences, Information Security Department
Imam Mohammad Ibn Saud Islamic University, Riyadh, Saudi Arabia*

e-mail: halnifisa@sm.imamu.edu.sa

As all the sectors, whether governmental or private are going for the digital transformation. And one of the most important sectors is the health sector. Consequently, the need for Electronic Medical Record (EMR) has increased significantly for the storage and analysis of medical data. This paper covers in depth the access controls for EMR specifically. The aim of this paper is to help institutions to select the proper access controls for them by comparing between the most popular existing access controls: Role Based Access Control (RBAC), Attribute Based Access Control (RBAC) and Blockchain. This comparison based on several criteria, which include dynamic environment, implementation, cloud, privacy and need to know, evaluated from the patient's privacy as well as from the institution's security requirements. The combination of the rated criteria and access controls produces the comparison matrix.

Keywords: EMR, Medical record, Access control, Security, Technology, RBAC, ABAC

POSTER PRESENTATIONS

PP1. FAST AND APPROXIMATE INFERENCE OF MULTILEVEL THRESHOLD AUTOREGRESSIVE MODEL FOR INTENSIVE LONGITUDINAL DATA VIA MEAN FIELD VARIATIONAL BAYES

Azizur Rahman^{1,2}, Depeng Jiang²

¹*Social Innovation Office, Department of Families, Government of Manitoba, Canada*

²*Department of Community Health Sciences, University of Manitoba, Winnipeg, Manitoba*

e-mail: emre.dunder@omu.edu.tr

The prevalence of diabetic disease is becoming one of the major health care problems in the world. To decrease the destroying effects of the diabetes, the glucose concentration should be kept under control. Several medications are used for this aim. We constructed a Bayesian lasso model for modeling the diabetic patient's health conditions. According to this model, we obtained the subset of medications which have significant effect on diabetes measurements. Bayesian lasso model also enables to find the diabetes condition of the patient's with determining the doses of the medications. It is shown that Bayesian lasso is more accurate than the classical lasso model.

Keywords: Bayesian approach, Lasso, Diabetic disease

PP2. ARTIFICIAL INTELLIGENCE-BASED MORPHOLOGY ANALYSIS SYSTEM FOR BRAIN ORGANOIDS

Elifsu Polatlı^{1,2}, Burak Kahveci^{1,2}, Sinan Güven^{1,2,3}

¹*Izmir Biomedicine and Genome Center, Dokuz Eylul University, 35340, Izmir, Türkiye*

²*Izmir International Biomedicine and Genome Institute, Dokuz Eylul University, 35340, Izmir, Türkiye*

³*Department of Medical Biology and Genetics, Faculty of Medicine, Dokuz Eylul University, 35340, Izmir, Türkiye*

e-mail: elifsu.polatli@ibg.edu.tr

Brain organoids are three-dimensional structures that mimic brain organogenesis and structure. Using brain organoids, neurodevelopmental and neurodegenerative disorders mechanisms can be studied, drug screening analysis can be performed. However, the experimental process of these organoids is quite expensive and laborious. It is very important to automate the process, which is open to human-based errors. In this direction, artificial intelligence-based methods can accelerate the process and minimize human-based errors. In this study, we developed the AI-based tool, which accelerates the brain organoid development process and can perform morphological analysis of organoids. This tool can classify organoids that can be used in the development process and those that are morphologically appropriate. In addition, it makes it possible to examine the steps in the development process by making morphological measurements with U-Net-based segmentation. It can detect specific structures such as neurons and rosette with object detection model. Automating many aspects of the brain organoid experimentation process, the tool provides a fast and comprehensive AI-based analysis.

Keywords: Computer vision, Deep learning, Brain organoids

PP3. TRANSCRIPTOMIC PROFILING OF INDUCED PLURIPOTENT STEM CELL DERIVED LACRIMAL ORGANOIDS

Burak Kahveci^{1,2}, Gamze Koçak^{1,2}, Canan Aslı Utine^{1,3}, Adil Mardinoğlu^{4,5},
Gökhan Karakülah^{1,2}, Sinan Güven^{1,2,6}

¹*Izmir Biomedicine and Genome Center, Izmir, Türkiye*

²*Izmir International Biomedicine and Genome Institute, Dokuz Eylul University, Izmir, Türkiye*

³*School of Medicine Department of Ophthalmology, Dokuz Eylul University, Izmir, Türkiye*

⁴*Science for Life Laboratory, KTH–Royal Institute of Technology, Stockholm, Sweden.*

⁵*Centre for Host-Microbiome Interactions, Faculty of Dentistry, Oral & Craniofacial Sciences, King's College London WC2R 2LS London, UK*

⁶*Department of Medical Biology School of Medicine, Dokuz Eylul University, Izmir, Türkiye*

e-mail: burak.kahveci@ibg.edu.tr

Ophthalmic diseases caused by visual impairments are a major problem affecting 2.2 billion people worldwide. Anterior segment of the eye is the most focused part in the eye research due to the critical role for vision quality. To study complex eye development processes, induced pluripotent stem cells (iPSCs) have potential to generate new approaches in regenerative medicine. The lacrimal gland is the main exocrine gland of the eye and secretes the lacrimal fluid, the main component of the eye film. This gland and its functions are quite important for eye health. Omics technologies have great potential to develop new differentiation protocols for anterior eye development to mimic native organs. In this study, we performed multi zonal ocular cell differentiation from human iPSCs for up to 45 days. Isolated mRNAs from lacrimal organoids at several time points of differentiation were analyzed at the transcriptome level and functionality of lacrimal organoids were assessed through stimulation with forskolin and carbachol. Our findings provide a comprehensive description of transcriptome landscape of multi zonal ocular cell differentiation for the first time.

Keywords: Transcriptomics, Induced pluripotent stem cells, Ophthalmology, Lacrimal organoids

PP4. PREDICTA: QUICK AND ACCURATE TRIAGE TOOL

Vahide Gül Türkmen¹, Burak Kahveci^{2,3}

¹*Faculty of Medicine, Erciyes University, 38030, Kayseri, Türkiye*

²*Izmir Biomedicine and Genome Center, Dokuz Eylul University, 35340, Izmir, Türkiye*

³*Izmir International Biomedicine and Genome Institute, Dokuz Eylul University, 35340, Izmir, Türkiye*

e-mail: vahidegulturkmen@gmail.com

The patient density in emergency services all over the world increases the workload of both doctors and healthcare professionals. The urgency levels of the patients who apply are different. The separation of these levels of urgency is called triage. Although triage distinguishes the urgency levels, the patient density is quite high since this process is done by people. Studies in which artificial intelligence methods perform health applications faster and more accurately are increasing day by day, and these methods can be used for triage in emergency services. In this study, we developed a machine learning-based system that performs triage faster and more accurately. We created this system using a synthetic dataset obtained with generative adversarial networks using real patient data. The artificial neural networks model developed with this dataset containing 41579 data achieved a very successful result with 98% recall and 99% AUC score.

Keywords: Machine learning, Synthetic data, Health data

PP5. MODIFIED CLINICAL KERNEL USING A COX MODEL

Seungyeoun Lee¹, Inyoung Kim², Hyunjae Lee¹

¹*Department of Mathematics and Statistics, Sejong University, Seoul, Korea*

²*Department of Statistics, Virginia Tech, USA*

e-mail: leesy@sejong.ac.kr

When using the support vector machines, traditional kernel functions such as the linear kernel are often applied to the set of the clinical data. However, these kernels do not consider the heterogeneous characteristics of clinical data, which is a mix of variable types with each variable its own range. A clinical kernel function was proposed by Daemen et al. (2009) to equalize the influence of clinical variables and take into account of the range of these variables. The clinical kernel function has shown to provide a better representation of patient's similarity and improve prediction of therapy response. In this paper, we propose a simple ensemble kernel by modifying the clinical kernel function using a Cox model. The proposed ensemble kernel differs from the clinical kernel function in that it takes a weighted average of the absolute effects of each clinical variable obtained after fitting the Cox model. We compare the performance of the ensemble kernel with that of the existing kernels on four data sets.

Keywords: Cox model, Clinical kernel, Ensemble kernel

PP6. NET BENEFIT IN CLINICAL DECISION MAKING PROCESS

Duygu Korkmaz Yalçın^{1,2}, İlker Ünal¹

¹*Çukurova University, Medicine Faculty, Department of Biostatistics*

²*Van Yüzüncü Yıl University, Medicine Faculty, Department of Medical Education and Informatics*

e-mail: duygukorkmaz@yyu.edu.tr

Decision Curve Analysis (DCA) is introduced by Vickers in 2006 to evaluate the clinical utility of diagnostic tests or treatment options. DCA is based on net benefit. The z-statistic is used to determine the statistical difference between the net benefit values calculated from 2 different models. In this study, the Friedman test was used instead of making a pairwise comparison each time for more than 2 models. For this, first of all, net benefit from each model (Generalized Linear Model, Decision Tree, Random Forest) was calculated using all the data, paired comparisons were made with the z test and the results were recorded. Then, different numbers of samples and repetitions were drawn from the data and the net benefit was calculated for each of the 3 models. Subsequently Friedman test was performed. Pima data (768 observations and 9 variables) were used as material. At a threshold probability of 0.01, the z-test found a statistical difference between at least 2 models. When this situation is taken as a reference, Friedman scenarios; observation number is n=100 and 50, 100, 200, 300 repetitions; n=200 and 50, 100, 200, 300 repetitions; n=300 and 50, 100 repetitions, are consistent with the Z test. For the first 102 observations and 9 variables, at a threshold probability of 0.01, there was no statistically difference between models as to z test. When this situation is taken as a reference, Friedman scenarios; n=30 and 10, 20, 30, 40, 50, 60, 70 repetitions; n=50 and 10, 20, 30, 40, 50, 60, 70 repetitions are consistent with the Z test. When there is a statistical difference, low number of observations and repetitions should be avoided in order to detect the difference. Also high number of observations and repetitions should be avoided in order to prevent finding a difference when there is no difference.

Keywords: Decision curve analysis, Model performance measure, Net benefit, Prediction models

PP7. ESTIMATION OF LOW-DENSITY LIPOPROTEIN CHOLESTEROL USING MACHINE LEARNING MODELS

Necla Koçhan¹

¹*Izmir Biomedicine and Genome Center*

e-mail: necla.kayaalp@gmail.com

Low-density lipoprotein cholesterol (LDL-C) is a commonly used measure in diagnosing and treating patients with dyslipidemia, which is frequently observed in children. Accurate estimation of LDL-C concentration is critical in diagnosing and treating cardiovascular disease (CVD), as LDL-C lowering therapy has been shown to reduce the risk of future coronary heart disease. Various formulas, such as Friedewald, Martin-Hopkins, Chen, Anandaraja, and Hattori, have been developed to estimate LDL-C values. However, many of these formulas have limitations and require further validation. Recently, machine learning (ML) methods such as random forest and support vector machines have been investigated for LDL-C estimation in different populations. Despite the application of ML approaches to LDL-C estimation, limited research has been conducted on the pediatric population, especially in Türkiye. Therefore, the aim of this study was to investigate the validity of LDL-C levels estimated by various LDL-C estimating formulas and powerful ML algorithms with directly measured LDL-C levels by Roche direct assay in the Turkish pediatric population. The study included 2,563 children under 18 years old who were treated at Sivas Cumhuriyet University Hospital in Sivas, Türkiye. The concordance between the estimates and direct measurements was assessed overall and separately for LDL-C and TG sublevels, and linear regression analyses were carried out, and residual error plots were created. The analysis showed that the ML models produced more concordant LDL-C estimates compared to LDL-C estimating formulas. In conclusion, accurate estimation of LDL-C concentration is crucial in CVD diagnosis and prognosis, and ML approaches such as random forest and support vector machines may provide more accurate LDL-C estimation in the pediatric population. Further research is required to validate these approaches in other populations and age groups.

Keywords: Cholesterol, Lipoproteins, Low-density lipoprotein, Triglyceride, Machine learning

PP8. EVALUATION OF SURVIVAL TREE RANDOM SURVIVAL FOREST AND COX PROPORTIONAL HAZARD MODELS

Duygu Korkmaz Yalçın¹, Sıddık Keskin¹

¹*Van Yüzüncü Yıl University, Medicine Faculty*

e-mail: duygukorkmaz@yyu.edu.tr

In the study, the performances of the Survival Tree and Random Survival Forest models that have some advantages (flexibility in assumptions) with Cox proportional hazards model were evaluated on the basis of the Concordance Index (CI). As application material, data sets in the R package program containing different variables and numbers of observations in the field of health; “cancer”, “veteran”, “diabetic”, “hd”, “swimsuit”, “GBSG2”, “cost”, “leukemia”, “kidney” were used. These datasets are randomly divided into 70% training and 30% test sets. In training datasets; The probability of survival for each observation on the test data was estimated using the Survival Tree, Random Survival Forest, and Cox proportional hazard models. The Concordance Index (CI) was calculated as a measure of concordance between estimated survival probabilities and actual survival times and event status (status). In order to increase the precision in the calculations, the average of the Concordance indices calculated after the data sets were randomly divided into the training and test sets 1000 times for all 3 models. "Rpart" and "Party" packages are used in the R package program for applications. The CI takes a value in the range of 0-1, and the closer it is to 1, the better the fit of the model. In tree-based models, the CI takes a maximum value of 0.5 when the splitting condition does not occur at all. Accordingly, for the Survival Tree model, the maximum and minimum CI was obtained from the “cancer” and “mayo” datasets respectively. Except for the “Kidney” data set, the highest CI value was obtained from the Survival Tree model in the other 8 data sets. In the Kidney data set, the highest CI value was obtained from the Random Survival Forest model. As a result, it was observed that the performance of the Survival Tree model was better in data sets with different number of variables and number of observations.

Keywords: Survival analysis, Survival tree, Random survival forest, Rpart

PP9. ALZHEIMER DISEASE CLASSIFICATION WITH MACHINE LEARNING METHOD

Fatma Gül Gezer¹, Berfu Parçalı², Kevser Setenay Öner², Fezan Mutlu²

¹*Istanbul University Cerrahpaşa, Medical Faculty, Biostatistics Department*

²*Eskişehir Osmangazi University, Medical Faculty, Biostatistics Department*

e-mail: gulgezer318@gmail.com

While the exact cause of Alzheimer's disease is not known so far, it is the most common cause of dementia. Prediction and early diagnosis of Alzheimer's disease; Although there is no known way to slow it down yet, it is very important in terms of the course of the disease and treatment methods. Classification of Alzheimer's Disease is very important in the treatment of symptoms and improving the quality of life of patients. Aim; In this study, machine learning algorithms and Alzheimer's data are classified. Clinical dementia rating was determined as the target variable, and 80% training and 20% test datasets were created. Support Vector Machine, Artificial Neural Networks and k-Nearest Neighbor algorithms were applied to the datasets. Classification was made with the applied algorithms and the results were compared. According to the comparison results, the classification performance is k-NN (AUC=0.575), SVM (AUC=0.612) and ANN (AUC=0.947), from lowest to highest, respectively. In the light of new studies on Alzheimer's Disease, new risk factors can be determined, clinical parameters are increased, studies can be expanded, and classification successes can be increased by transferring them to machine learning algorithms.

Keywords: Alzheimer disease, Machine learning, Support vector machine, Artificial neural network, k-Nearest neighbor

PP10. A COMPREHENSIVE R SHINY WEB TOOL THAT COMBINES TWO CONTINUOUS DIAGNOSTIC TESTS

Serra İlayda Yerlitaş^{1,2}, Serra Bersan Gengeç², Gözde Ertürk Zararsız^{1,2}
Selçuk Korkmaz³, Gökmen Zararsız^{1,2}

¹*Department of Biostatistics, Faculty of Medicine, Erciyes University, Kayseri, Türkiye*
²*Drug Application and Research Center (ERFARMA), Erciyes University, Kayseri, Türkiye*
³*Department of Biostatistics, Faculty of Medicine, Trakya University, Edirne, Türkiye*

e-mail: serrayerlitas@erciyes.edu.tr

Diagnostic tests are widely used in the health sector to diagnose diseases. The diagnostic accuracy, performance and reliability of these tests are considered when making diagnostic tests widely available. For diseases with more than one diagnostic test, the most common approach is to compare the performance of the diagnostic tests or ratio the diagnostic tests to try to improve performance. Highly accurate estimates can be obtained by combining diagnostic tests using statistical methods. Although there are many methods in the literature for combining diagnostic tests, there is no comprehensive web software that can apply these methods. The aim of this study is to develop an easy-to-use, free, and comprehensive shiny web application using the dtComb R library developed for combining two diagnostic tests. The dtComb, shiny, stringr, shinyBS, dplyr, and ggplot2 libraries of the R programming language were used to develop the web application, and the shinytest library was used for the testing process. This developed application has been transferred to the workstation in Erciyes University Faculty of Medicine, Department of Biostatistics, and opened for use. The dtComb web software allows users to perform calculations quickly and easily using a data set containing the values of the gold standard and two biomarkers as input. Thanks to this software, users can quickly combine diagnostic tests with a total of 143 combination methods under 4 main headings: 8 linear combination methods, 7 nonlinear combination methods, 14 mathematical operator methods, and 113 machine-learning algorithms. In addition, it can enhance the performance of diagnostic test statistics with 5 standardization methods and 12 resampling methods prior to calculation. Users can report the results in various file formats and present them with comparative graphs. The dtComb web application can be accessed at <http://biosoft.erciyes.edu.tr/app/dtComb>.

Keywords: Combining diagnostic tests, Machine learning, ROC analysis, R programming language

PP11. A CIRCULAR HEATMAP VISUALIZATION APPROACH FOR INTERLABORATORY COMPARISONS IN RING STUDIES

Gözde Ertürk Zararsız^{1,2}, Alexander Cecil³, Jutta Lintelmann², Gernot Poschet⁴, Jennifer Kirwan⁵; Sven Schuchardt⁶, Xue Li Guan⁷, Daisuke Saigusa⁸, David Wishart⁹, Jiamin Zheng¹⁰, Rupasri Mandal¹⁰; Lisa St. John-Williams¹¹, Kendra Adams¹¹, J. Will Thompson¹¹, Michael P. Snyder¹², Kevin Contrepois¹², Songlie Chen¹², Nadia Ashrafi¹³, Sumeyya Akyol¹³, Ali Yilmaz¹³, Stewart Graham¹³, Thomas M. O'Connell¹⁴, Karl Kalecky^{15,16}, Teodoro Bottiglieri^{15,1}, Tuan Hai Pham¹⁷; Therese Koal¹⁷, Jerzy Adamski^{18,19,20}, Gabi Kastenmüller²

¹Department of Biostatistics, Erciyes University School of Medicine, 38039 Kayseri, Türkiye, Drug Application and Research Center (ERFARMA), Erciyes University, 38280 Kayseri, Türkiye, Institute of Computational Biology, Helmholtz Zentrum München - German Research Center for Environmental Health, Ingolstädter Landstraße 1, 85764 Neuherberg, Germany

³Metabolomics and Proteomics Core, Helmholtz Zentrum München, German Research Center for Environmental Health, Ingolstädter Landstraße 1, 85764 Neuherberg, Germany

⁴Metabolomics Core Technology Platform, Heidelberg University, Germany

⁵Berlin Institute of Health, Metabolomics Platform, Germany

⁶Fraunhofer-Institute for Toxicology and Experimental Medicine, Hannover, Germany

⁷Lee Kong Chian School of Medicine, Nanyang Technological University, Singapore

⁸Laboratory of Biomedical and Analytical Sciences, Faculty of Pharma-Science, Teikyo University, Japan

⁹Department of Biological Sciences, University of Alberta, Alberta, Canada

¹⁰Department of Computing Sciences, University of Alberta, Alberta, Canada

¹¹Duke University, Durham, NC, USA

¹²Stanford University, Stanford, CA, USA

¹³Beaumont Health System, Royal Oak, MI, USA

¹⁴Indiana University, Indianapolis, IN, USA

¹⁵Center of Metabolomics, Institute of Metabolic Disease, Baylor Scott & White Research Institute, Dallas, TX, USA

¹⁶Institute of Biomedical Studies, Baylor University, Waco, TX, USA

¹⁷biocrates life sciences ag, Innsbruck, Austria

¹⁸Institute of Experimental Genetics, Helmholtz Zentrum München, German Research Center for Environmental Health, Ingolstädter Landstraße 1, 85764 Neuherberg, Germany

¹⁹Department of Biochemistry, Yong Loo Lin School of Medicine, National University of Singapore, 8 Medical Drive, Singapore 117597, Singapore

²⁰Institute of Biochemistry, Faculty of Medicine, University of Ljubljana, Vrazov trg 2, 1000 Ljubljana, Slovenia

e-mail: gozdeerturk@erciyes.edu.tr

Many statistical approaches are used and developed for the assessment of analytical variation between international laboratories. Presentation and visualization of measurement errors with statistical models is of great importance in evaluating the analytical performance of measurements. The aim of this study is to introduce a graphical approach used to compare the performance of measured values in international laboratories between laboratories and by markers. Real data is a data set created with the aim of evaluating the analytical performance of the MxP® Quant 500 (Biocrates Life Science AG, Innsbruck, Austria) metabolomics kit, which is used in the quantification of nearly 634 metabolites in 14 international laboratories. Each laboratory participating in the study was given a single letter label from A to N. Data were collected from the laboratories where the experiments were carried out.

Interlaboratory reproducibility was evaluated with coefficients of variation. The difference between each laboratory was evaluated by method comparison. Analytical measurement error rates were assumed to be constant between laboratories. Deming regression method was used for comparisons between laboratory pairs. Measurement errors were estimated from triplicate measurements of each metabolite. The jackknife method was used to determine the confidence intervals. Absolute and relative differences were calculated with linear regression analyses. The results were evaluated separately for the minimum, median and maximum values of each metabolite. These results are visualized with circular heatmap graphs. The researcher was given a graphical report on which metabolite level between laboratories and which laboratory measured with the highest error.

Keywords: Relative difference, Deming regression, Metabolomics, Systematic variation, Method comparison

PP12. A META-ANALYSIS STUDY FOR DIAGNOSING SKIN CANCER WITH MACHINE LEARNING TECHNIQUES

Gözde Ertürk Zararsız^{1,2}, Elif Çelik Gürbulak^{2,3}, Serra İlayda Yerlitaş^{1,2}, Selen Yılmaz Işıktan⁴, Abdullah Demirbaş⁵, İrem Eroğlu^{2,6}, Aleyna Erakçaoğlu^{2,6}, Ragıp Ertaş⁷, Ömer Faruk Elmas⁸, Gökmen Zararsız^{1,2}

¹ Department of Biostatistics, Erciyes University School of Medicine, Kayseri, Türkiye,

² Drug Application and Research Center (ERFARMA), Erciyes University, Kayseri, Türkiye

³ Erciyes University, Faculty of Veterinary Medicine, Department of Biometrics, Kayseri, Türkiye

⁴ Vocational Higher School of Social Sciences, Hacettepe University, Ankara, Türkiye; ⁴ Department of Biostatistics, Faculty of Medicine, Hacettepe University, Ankara, Türkiye

⁵ Selçuk University, Faculty of Medicine, Department of Dermatology, Konya, Türkiye

⁶ Faculty of Life and Natural Sciences Department of Molecular Biology and Genetics, Abdullah Gul University, Kayseri, Türkiye

⁷ Kayseri Education and Research Hospital Dermatology Clinic, Kayseri, Türkiye

⁸ Kırşehir Ahi Evran University, Department of Dermatology, Kırşehir, Türkiye

e-mail: gozdeerturk@erciyes.edu.tr

Artificial intelligence is a technology which uses machines and programs to simulate human behavior. Supervised learning is the most common type of learning used in dermatology. In this study, SCOPUS was performed with 1888 primary studies searched in the MEDLINE and Web of Science databases with the search strategy in Appendix 1. Inclusion criteria; The use of dermatoscopy or dermoscopy methods in diagnosis, the use of AI algorithms, the subject of the study being melanoma, the language of publication in English, the use of skin lesion images in the study, the classification. By removing studies that did not meet the inclusion criteria from the study sample, the meta-analysis consisted of 79 diagnostic test data from 35 primary studies. Moderator variables used in the study; segmentation was determined as analytical models. The quality assessment of the studies was done with the QUADAS-2 tool. Bivariate meta-analysis technique was used to obtain common sensitivity, common specificity, joint positive odds ratio for skin cancer detection by machine learning. These estimates are graphically represented by the forest plot (sensitivity and specificity). Summary Receiver Operator Characteristics curves (sROC) were also created to understand the diagnostic value of skin cancer detection by machine learning. Fagan chart was created to find out whether the result of skin cancer detection by machine learning changes the probability of skin cancer in the suspected patient. Heterogeneity was defined statistically using the Chi-square-based Q statistic. I² statistics were calculated to measure the level of inconsistency. Significant heterogeneity was detected between studies. According to the results of the meta-analysis, the common sensitivity rate is 0.89 [95% CI: 0.87 - 0.91], the common specificity rate is 0.94 [95% CI: 0.92, 0.95], the common positive probability rate is 14.5 [95% CI: 11.1 - 18.9], the common the negative odds ratio was 0.12 [95% CI: 0.10, 0.14] and the diagnostic odds ratio was 128 [95% CI: 83 - 197] and the ROC area was 0.97 [95% CI: 0.94 - 0.98].

Keywords: Melanoma, Dermatology, Bivariate meta-analysis, Machine learning techniques

PP13. ESTABLISHING CONTINUOUS REFERENCE INTERVALS FOR THYROID FUNCTION TESTS

Funda İpekten^{1,2,3}, Gözde Ertürk Zararsız^{1,2}, Halef Okan Doğan⁴, Çiğdem Karakükçü⁵, Gökmen Zararsız^{1,2}

¹*Department of Biostatistics, Faculty of Medicine, Erciyes University, Kayseri, Türkiye*

²*Drug Application and Research Center, Erciyes University, Kayseri, Türkiye*

³*Department of Biostatistics, Faculty of Medicine, Adiyaman University, Adiyaman, Türkiye*

⁴*Department of Biochemistry, Faculty of Medicine, Cumhuriyet University, Sivas, Türkiye*

⁵*Department of Biochemistry, Faculty of Medicine, Erciyes University, Kayseri, Türkiye*

e-mail: fundaipekten@gmail.com

Clinical laboratory tests are an important part of the healthcare field. Interpretation of clinical laboratory test values is dependent on the availability of accurate health-related measures, known as reference intervals. In determining the reference intervals and interpreting the tests, basic variables such as age and gender are especially considered. Reference intervals may not give complete and accurate results due to variations that may occur due to gender and age in clinical laboratory tests. In line with these results, wrong treatment methods can be applied or wrong diagnoses can be made. Therefore, reference intervals should be studied separately according to basic variables such as age and gender. The complex relationship between clinical laboratory tests and age can be made simple with reference intervals (Hall et al., 2021). Continuous reference intervals should be used to represent this relationship more accurately. In our study, we found this relationship in thyroid function tests (Thyroid-stimulating hormone (TSH), free thyroxine (fT4) We aimed to establish continuous reference intervals for free triiodothyronine (fT3) and compare the performances of these methods. The study used thyroid function test data from the pediatric group in the same analytical system between 2017 and 2022. Methods BCPE (μ, σ, ν, τ), BCT (μ, σ, ν) and BCCG (μ, σ, ν) were used to construct continuous reference intervals. The generalized Akaike information criterion (GAIC) was used in the performance comparisons of the methods. In establishing the continuous reference intervals for the thyroid function tests used in our study, the performance success of the BCT method was further evaluated. It can be said that it is comparatively better than other methods.

Keywords: GAMLSS, Continuous reference intervals, Precision, Thyroid function tests

PP14. ON PERFORMANCES OF DIFFERENT CORRELATION COEFFICIENTS

Yeşim Uzun Uğur¹, Mehmet Mendeş¹

¹*Çanakkale onsekiz Mart University, Agriculture Faculty, Biometry and Genetics Unit,
17100, Çanakkale/ Türkiye*

e-mail: yesimuu@gmail.com

Since there are many different correlation coefficients have been developed and proposed for different cases, it is extremely important to aware of which correlation coefficient(s) is more appropriate. This study aimed at investigating the performances of eight different correlation coefficients namely Pearson, Spearman's Rank, Percentage Bend, Winsorized, Distance, Blomqvist, Biweight Midcorrelation, Hoeffding's D under different experimental conditions via a comprehensive Monte Carlo Simulation Study. Results of the simulation study showed that the performances of these correlation coefficients are affected by sample size, effect size, and distribution shape. When both the type I error and test power estimates are evaluated together, it is seen that the Pearson's, Winsorized, and Spearman Rank correlation coefficients are generally the most appropriate coefficients for many experimental conditions. However, as the sample size increases all coefficients except for Blomqvist tend to give very similar estimates. The most deviated estimates, on the other hand, have been obtained when Bloqvist correlation is used.

Keywords: Correlation, Type I error rate, Test power, Simulation

YOUNG STATISTICIANS SHOWCASE PRESENTATIONS

MODIFIED SHAP METHOD FOR SEASONAL VACCINATION STATUS

Ahmet Yalcin¹, Bekir Cetintav², Selim Cetin³

¹*Burdur Mehmet Akif Ersoy University, The Graduate School of Natural and Applied Sciences*

²*Burdur Mehmet Akif Ersoy University, Department of Statistics*

³*Burdur Mehmet Akif Ersoy University, Department of Mathematics*

e-mail: ahmtylcinn15@gmail.com

Seasonal vaccination is important for people at risk of severe illness from the influenza/flu. Every year, flu vaccination prevents illnesses, medical visits, hospitalizations, and deaths. This is considered an important problem that affects the success of vaccination programs. In the literature, several studies have examined the desire to vaccinate against swine flu in different countries. In these studies, the behavioral data obtained mostly by surveys/questionnaires were analyzed with statistical and machine learning (ML) methods. The results obtained draw a global framework about factors affecting the willingness of vaccination. In the light of these data, more “general” approaches are developed in the field of public health. In our study, the development of “personalized” approaches are discussed. Factors affecting the target variable (in our study, target variable is “get a seasonal vaccination”) in ML models are measured in two ways: feature importance methods for global interpretations, “explainable artificial intelligence (XAI) methods” for local interpretations. In our study, we set up a ML model for estimation whether a person will get a seasonal vaccine. In addition, by applying XAI methods, we can predict the factors that affect each person’s condition. In this way, “personalized” recommendations that will increase the likelihood of people to get a seasonal vaccine can be produced, and new preventive medicine applications can be developed. We used SHAP, an XAI method, but it has some limitations/disadvantages discussed in the literature. In this work, we propose a new approach for the case of under-interpretation in ordered classes, which is another disadvantage of SHAP for the seasonal vaccination case that has not been discussed in the literature: We construct a new linear model for the additive feature attribution method, while keeping the calculation of Shapley values as proposed by Lundberg and Lee.

Keywords: Machine Learning, XAI (explainable artificial intelligence), SHAP, Vaccination

SHINY APP FOR GO ANALYSIS (SimElegans)

İrem Kahveci¹, Hamdi Furkan Kepenek¹, Dinçer Göksülük²

¹*Abdullah Gül University Molecular Biology and Genetics*

²*Erciyes University School of Medicine*

e-mail: irem.kahveci@biogenr.com

The use of *C. elegans* as a model organism in human-oriented studies is of great importance. *Caenorhabditis elegans* is a useful model organism for studying a variety of diseases ranging from SMA disease to mitochondrial diseases. It contributes to science in the analysis of data obtained in image processing research. R Shiny package, which is a kind of open software and increasingly used in data science, is a powerful tool that enables to create interactive web applications directly from R with results in various research fields such as genetics, engineering, bioinformatics. The main purpose of this study is to visualize the common genes between human and *C. elegans* organism on the 3D *C. elegans* image for their similarities in terms of ontology and localization.

Keywords: R Shiny, GO Similarity, *C. elegans*

REAL-TIME DETECTION OF THE START AND SUBSEQUENT EPIDEMIC STATES OF HIV OUTBREAKS AMONG PEOPLE WHO INJECT DRUGS: INSIGHTS FROM FOUR EUROPEAN COUNTRIES

Valia Baralou^{1*}, Argiro Karakosta^{1*}, Christos Thomadakis¹, Nikos Demiris², Nikos Pantazis¹, Olga Anagnostou³, Christos Danopoulos³, Dimitris Katsiris⁴, Giota Touloumi¹

¹*Department of Hygiene, Epidemiology & Medical Statistics, Medical School, National & Kapodistrian University of Athens, Athens, Greece*

²*Department of Statistics, Athens University of Economics and Business, Athens, Greece*

³*Greek Organisation Against Drugs (OKANA), Athens, Greece*

⁴*InDigital S.A., Athens, Greece*

*These authors contributed equally to this work.

e-mail: vbaralou@med.uoa.gr

HIV outbreaks among people who inject drugs are not uncommon, but differ in their magnitude, duration and post-epidemic level. We aimed to compare the performance of different methods for real-time detection of growth, decline and post-epidemic states of such outbreaks in four European countries. Data on weekly HIV diagnoses in Greece, Romania, Luxembourg and Ireland were provided by The European Surveillance System of the European Centre for Disease Prevention and Control. To identify each state, we developed a two-state hidden Markov model (HMM), a method based on prediction interval (PI) and a novel approach that combines the former two by assigning weights to HMM's transition probabilities at each time point based on the PI limits (HMM-PI). We also applied control charts for anomaly detection. Methods were applied prospectively; performance was assessed with sensitivity, specificity and timeliness (i.e., the difference between the start of each state and the first alarm after its onset). Concerning the extremely small outbreaks in Luxembourg and Ireland that returned to pre-epidemic levels, only control charts identified the growth and post-epidemic state. Contrarily, when applied for detecting the large and abrupt outbreaks in Greece and Romania, PI performed at least similarly with classic methods. HMM-PI had also excellent performance albeit with moderate specificity of detecting the Romania outbreak. HMM-PI gave the best balance among all metrics for the decline phase detection. For detecting the post-epidemic state, no method reached a satisfying balance between these metrics but PI performed better giving excellent sensitivity and timeliness but moderate specificity. No method is a panacea for identifying all states of an outbreak. The proposed methods HMM-PI and PI seem promising for monitoring large and explosive epidemics that do not reach the pre-epidemic level while classic control charts are adequate for the surveillance of very small outbreaks.

Keywords: Anomaly detection, HIV surveillance, Hidden Markov model, Control charts

MODEL BASED CLUSTERING FOR SPATIAL DATA

Anna Nalpantidi¹

¹*Department of Statistics, Athens University of Economics and Business,
76 Patision St., 10434, Athens, Greece*

e-mail: *analpa@aueb.gr*

Clustering spatial data using standard finite mixture models is not always an efficient way. Spatial heterogeneity and spatial dependence, the two main characteristic of spatial data have to be taken into consideration. Most of the proposed models in literature refer to poisson or normal mixtures. The purpose of this paper is to extend standard finite mixture models in the context of multinomial mixture for spatial data, in order to cluster geographical units according to demographic characteristics. The spatial information is incorporated on the model through mixing probabilities of each component. To be more specific, a Gibbs distribution is assumed for prior probabilities. In this way, assignment of each observation is let to be affected by neighbors' cluster and spatial dependence is included in the model. Estimation is based on a modified EM algorithm which is enriched by an extra, initial step for approximating the field. The simulated field algorithm is used in this initial step. The presented model will be used for clustering municipalities of Attica with respect to age distribution of residents.

Keywords: Model based clustering, Multinomial mixture, EM algorithm, Markov random fields

MULTILEVEL BAYESIAN NETWORK TO MODEL CHILD MORBIDITY USING GIBBS SAMPLING

Bezalem Eshetu Yirdaw¹, Legesse Kassa Debusho¹

¹*University of South Africa, c/o Christiaan de Wet Road & Pioneer Avenue, Johannesburg, South Africa*

e-mail: 12962805@gmail.com

Child morbidity has been a serious global burden, especially in sub-Saharan African countries. Bayesian networks (BNs) are suitable models for studying multiple outcomes, such as more than one child morbidities, simultaneously. However, these models fail the assumption of independent observation in the case of hierarchical data. Therefore, this study proposes a random intercept multilevel Bayesian network (MBN) models to study the conditional dependencies between multiple outcomes. The structure of MBN was learned using the connected three parent set block Gibbs sampler, where each local network was included based on Bayesian information criteria (BIC) score of multilevel regression. The model was examined using simulated data assuming features of both multilevel models and BNs. The estimated area under the receiver operating characteristics for both models were above 0.8, indicating good fit. The MBN was then applied to real child morbidity data from the 2016 Ethiopian Demographic Health Survey (EDHS). The result shows a complex causal dependency between malnutrition indicators and child morbidities such as anemia, acute respiratory infection (ARI) and diarrhea. According to this result, families and health professionals should give special attention to children who suffer from malnutrition and also have one of these illnesses, as the co-occurrence of both can worsen a child's health.

Keywords: Bayesian network, Child morbidity, Directed acyclic graph, Multilevel Bayesian network, The Connected three parent Gibbs sampling

INDEX

- A. Ergun Karaagaoglu, 4, 5, 71
 Abdullah Demirbaş, 94
 Adil Mardinoğlu, 84
 Ahmet Mert Saygu, 78
 Ahmet Ömer Özgür, 78
 Ahmet Öztürk, 3, 43
 Ahmet Sezgin, 53
 Ahmet Yalcin, 98
 Ahu Cephe, 53, 76
 Aldo Cocco, 23
 Alexander Cecil, 92
 Aleyna Erakçaoğlu, 94
 Ali Yilmaz, 92
 Almond Stöcker, 5, 12
 Amit Huppert, 9, 41, 77
 Anja Victor, 72
 Anna Nalpantidi, 101
 Anne-Laure Boulesteix, 5, 6, 49
 Argiro Karakosta, 100
 Aris Perperoglou, 5, 13
 Arne Bathke, 9
 Arrigo F.G. Cicero, 76
 Atilla H. Elhan, 3, 43
 Ayca Olmez, 54
 Aylin Alın, 54, 59, 78
 Azizur Rahman, 82
 Bahar Taşdelen, 3
 Bahjat F. Qaqish, 64
 Banu Isbilen Basok, 68
 Bekir Cetintav, 98
 Bella Vakulenko-Lagun, 3
 Benjamin Reiser, 3, 5, 32
 Berfu Parçalı, 90
 Berk Pişkin, 78
 Bezalel Eshetu Yirdaw, 102
 Bhramar Mukherjee, 9
 Birol Emir, 3
 Buğra Varol, 66
 Buket İpek Berk, 74
 Burak Kahveci, 83, 84, 85
 Burak Kürsad Günhan, 72
 Burcu Bakir-Gungor, 64, 65
 Bülent Yılmaz, 58
 Büşra Emir, 73
 Cagdas Hakan Aladag, 5, 14
 Canan Aslı Utine, 84
 Cemil Çolak, 3
 Cengiz Bal, 3
 Chanhee Lee, 67
 Christos Danopoulos, 100
 Christos T. Nakas, 3
 Christos Thomadakis, 100
 Constantine Gatsonis, 3
 Cuneit Ozden, 61
 Çiğdem Karakükçü, 95
 Daisuke Saigusa, 92
 David M. Steinberg, 5, 9, 15
 David Refaeli, 15
 David Wishart, 92
 David Zucker, 3
 Dean Palejev, 45
 Demet Arı, 74
 Denitsa Grigorova, 45
 Deniz Ilhan Topcu, 8, 61, 68
 Depeng Jiang, 82
 Didem Turgut, 61
 Dimitar Vassilev, 62
 Dimitris Karlis, 3, 51, 52
 Dimitris Katsiris, 100
 Dinçer Göksülük, 4, 55, 71, 99
 Douglas Midthune, 30
 Duygu Korkmaz Yalçın, 60, 87, 89
 Ebru Aker, 58
 Ebru Kaya Başar, 48
 Elena Nardi, 79
 Elif Çelik Gürbulak, 94
 Elif Kaymaz, 73
 Elifsu Polatlı, 83
 Ella Shaposhnik, 15
 Emrah Gecili, 5, 16
 Erdal Coşgun, 4, 8
 Erdem Karabulut, 3, 53
 Esmâ Gamze Aksel, 57
 Esra Kutsal Mergen, 42
 Fatimah M. Alturkistani, 80
 Fatma Ezgi Can, 73
 Fatma Gül Gezer, 90
 Ferhan Elmalı, 3, 73
 Fezan Mutlu, 90
 Figen Demirkazık, 69
 Fikret Er, 3
 Filomena Maggino, 9
 Funda İpekten, 95
 Gabi Kastenmüller, 92
 Gamze Durhan, 69
 Gamze Koçak, 84
 Geert Molenberghs, 3, 5, 9, 17, 33
 Georgia Rompoti, 51
 Gernot Poschet, 92
 Gilad Shapira, 15
 Giorgos Bakoyannis, 5, 18
 Giota Touloumi, 100
 Gordon Smyth, 5, 7
 Gökhan Karakülâh, 84
 Gökmen Zararsız, I, 4, 43, 53, 57, 76, 91, 94, 95
 Göknuur Giner, 4, 5, 19

- Gözde Ertürk Zararsız, 4, 43, 53, 76, 91, 92, 94, 95
Grant Izmirlian, 30
Guadalupe Gómez Melis, 5, 4
Günel Bilek, 70
H. Refik Burgut, 4, 5
Halef Okan Doğan, 76, 95
Hamdi Furkan Kepenek, 55, 99
Hamparsum Bozdogan, 5, 20, 34
Hanife Avcı, 69
Havi Murad, 3, 5, 21
Hernando Ombao, 3, 29
Hilah Alnafisah, 80
Hyunjae Lee, 86
Inyoung Kim, 86
Itai Dattner, 3
Ivor Cribben, 26
İhsan Berk, 74
İlker Ercan, 3
İlker Ünal, 4, 87
İmran Kurt Ömürlü, 66
İrem Eroğlu, 94
İrem Kahveci, 55, 99
J. Will Thompson, 92
Jale Karakaya, 69
Jaroslaw Harezlak, 3
Jay Kaufman, 49
Jennifer Kirwan, 92
Jerzy Adamski, 92
Jiamin Zheng, 92
Jutta Lintelmann, 92
Kaloyan Vitanov, 46
Karen Kafadar, 5, 8
Karl Kalecký, 92
Kenan Köse, 3
Kendra Adams, 92
Kevin Contrepolis, 92
Kevser Setenay Öner, 90
Konstantinos Fokianos, 3, 4
KyungMann Kim, 5, 23
Legesse Kassa Debushe, 47, 102
Leta Lencha Gemechu, 47
Lisa St. John-Williams, 92
Lisa Steyer, 12
Lu Mao, 23
M.Yasemin Akşehirli Seyfeli, 43
Malgorzata Bogdan, 5, 24
Malik Yousef, 64, 65
Maria-Tereza Dellaporta, 52
Marie Trussart, 7
Marie-Eve Beauchamp, 49
Marina Bogomolov, 44
Maroussia Slavtchova-Bojkova, 46
Maya Zhelyazkova, 62
Mehmet Ali Kaygusuz, 50, 63
Mehmet Gönen, 5, 25
Mehmet Koçak, 5, 35
Mehmet Mendeş, 96
Mehmet Orman, 3
Meltem Gülsün Akpınar, 69
Meng Wu, 40
Meriç Yavuz Çolak, 3
Merve Basol Goksuluk, 71
Merve Kasikci, 4, 56
Mevlüt Türe, 66
Micha Mandel, 77
Michael P. Snyder, 92
Michal Abrahamowicz, 49
Mira Marcus-Kalish, 15
Mithat Gönen, 3
Mustafa Agah Tekindal, 48, 73
Mustafa Çavuş, 4
Naama M. Kopelman, 41, 77
Nadia Ashrafi, 92
Necla Koçhan, 4, 53, 76, 88
Nikos Demiris, 100
Nikos Pantazis, 100
Nirit Agay, 21
Nurten Bulut, 64
Olga Anagnostou, 100
Orhun Öztürk, 43
Ori Davidov, 3
Osman Dag, 4, 56
Ozgur Saman, 56
Ozlem İlk, 36
Ömer Faruk Elmas, 94
Özlem İlk, 3, 5
Pavel Mozgunov, 72
Philip Tzvi Reiss, 3, 5, 26
Pınar Günel, 74
Pınar Özdemir, 3
Ping Hu, 5, 9
Przemyslaw Biecek, 6
R. Todd Ogden, 5, 27
Rachel Dankner, 21
Ragıp Ertay, 94
Raina E. Josberger, 40
Ralitza Gueorguieva, 45
Ralph Brinks, 9
Ramyar Molania, 5, 7, 28
Rawaby Alsaaid, 80
Recai M. Yücel, 5, 37, 40
Refika Sultan Doğan, 58
Rhonda Szczesniak, 16
Rim Khazhin, 78
Rupasri Mandal, 92
Ruth Heller, 3, 44
Samet Senel, 61
Selçuk Korkmaz, 91

- Selen Yılmaz Işıkkhan, 94
Selim Can Kuralay, 57
Selim Cetin, 98
Serkan Bolat, 76
Serkan Doğan, 58
Serra Bersan Gengeç, 91
Serra İlayda Yerlitaş, 76, 91, 94
Seungyeoun Lee, 86
Sevilay Karahan, 4, 42
Sharon Amit, 41, 77
Sıddık Keskin, 3, 60, 89
Sinan Güven, 83, 84
Sipan Aslan, 5, 29
Songlie Chen, 92
Sonja Greven, 12
Stefan Tsonev, 62
Stewart Graham, 92
Sumeyya Akyol, 92
Sven Schuchardt, 92
Taesung Park, 67
Teodoro Bottiglieri, 92
Therese Koal, 92
Thomas M. O'Connell, 92
Tim Morris, 49
Tuan Hai Pham, 92
Tuğba Özkal Yıldız, 75
Tuo Wang, 23
Urania Dafni, 3, 51, 52
Ünal Erkorkmaz, 3
Vahap Eldem, 57
Vahide Gül Türkmen, 85
Valia Baralou, 100
Victor Kipnis, 5, 30
Vilda Purutçuoğlu, 4, 50, 63
Vildan Sümbüloğlu, 3, 74
Vincent Carey, 5, 8, 2
Willi Sauerbrei, 49
Xu Meng, 26
Xue Li Guan, 92
Yahya Laleli, 5, 38
Yair Goldberg, 77
Yasin Görmez, 60
Yavuz Sanisoğlu, 3
Yeşim Uzun Uğur, 96
Yoav Benjamini, 3, 5, 10
Yonatan Woodbridge, 41, 77
Yusuf Can Sevil, 75
Zeynep Baykan, 43
Zeynep Özel, 48



EMR2023
İZMİR, TÜRKİYE | MAY 8-11, 2023



etix
events