

Student Contributed Presentations

Oral presentations:

1. S-GMAS: Shape based Genome-wide Mediation Analysis

Shengxian Ding^{1*}, Rongjie Liu¹, Anuj Srivastava¹ and Chao Huang¹

¹Department of Statistics, Florida State University

Causal mediation analysis is widely utilized in neuroscience to investigate the role of brain image phenotypes in the neurological pathways from genetic exposures to clinical outcomes. However, it is still difficult to conduct a genome-wide mediation analysis with the shapes of brain regions as mediators due to several challenges including (i) large-scale genetic exposures, i.e., millions of single-nucleotide polymorphisms (SNPs); (ii) nonlinear Hilbert space for shape mediators; and (iii) statistical inference on the direct and indirect effects. To tackle these challenges, this paper proposes a mediation analysis framework with high dimensional genetic exposures and shape mediators. First, the square-root velocity function representations are extracted from the shapes, which fall in an unconstrained linear Hilbert subspace. Second, to address the issue caused by the high dimensionality in genetic exposures, the global sure independence screening procedure is conducted to discover candidate SNPs influencing the shape mediators. To identify the underlying causal pathways from the detected SNPs to the clinical outcome implicitly through the shape mediators, we proposed a framework consisting of a shape-on-scalar model and a scalar-on-shape model. Furthermore, the bootstrap resampling approach is adopted to investigate both global and local significant mediation effects. Finally, our framework is applied to the ADNI corpus callosum shape data and we successfully identify the mediation effect of a subset of candidate SNPs on Alzheimer's Disease through a subregion of the corpus callosum.

2. A Novel Causal Mediation Analysis Approach for Zero-Inflated Mediators

Meilin Jiang, Department of Biostatistics, University of Florida

Mediation analyses play important roles in making causal inference in biomedical research to examine causal pathways that may be mediated by one or more intermediate variables (i.e., mediators). Although mediation frameworks have been well established such as counterfactual-outcomes (i.e., potential-outcomes) models and traditional linear mediation models, little effort has been devoted to dealing with mediators with zero-inflated structures due to challenges associated with excessive zeros. We develop a novel mediation modeling approach to address zero-inflated mediators containing true zeros and false zeros. The new approach can decompose the total mediation effect into two components induced by zero-inflated structures: the first component is attributable to the change in the mediator on its numerical scale which is a sum of two causal pathways and the second component is attributable only to its binary change from zero to a non-zero status. An extensive simulation study is conducted to assess the performance and it shows that the proposed approach outperforms existing

standard causal mediation analysis approaches. We also showcase the application of the proposed approach to a real study in comparison with a standard causal mediation analysis approach.

3. JUMP: replicability analysis of high-throughput experiments with applications to spatial transcriptomic studies

Pengfei Lyu^{1*}, Yan Li², Xiaoquan Wen³, Hongyuan Cao¹

¹Department of Statistics, Florida State University

²Jilin University

³University of Michigan

Replicability is the cornerstone of scientific research. The current statistical method for high-dimensional replicability analysis either cannot control the false discovery rate (FDR) or is too conservative. We propose a statistical method, JUMP, for the high-dimensional replicability analysis of two studies. The input is a high-dimensional paired sequence of p-values from two studies and the test statistic is the maximum of p-values of the pair. JUMP uses four states of the p-value pairs to indicate whether they are null or non-null. Conditional on the hidden states, JUMP computes the cumulative distribution function of the maximum of p-values for each state to conservatively approximate the probability of rejection under the composite null of replicability. JUMP estimates unknown parameters and uses a step-up procedure to control FDR. By incorporating different states of composite null, JUMP achieves a substantial power gain over existing methods while controlling the FDR. Analyzing two pairs of spatially resolved transcriptomic datasets, JUMP makes biological discoveries that otherwise cannot be obtained by using existing methods.

4. Death or Survival: Which You Measure May Affect Conclusions

Jake Shannin^{1*}, Babette A. Brumback²

¹Department of Statistics, University of Florida

²Department of Biostatistics, University of Florida

Background and Aims: Considering the opposite outcome, for example, survival instead of death, may affect conclusions about which subpopulation benefits more from a treatment or suffers more from an exposure. Methods: For a case study on bankruptcy following melanoma, we compute and interpret the relative risk, odds ratio, and risk difference for different age groups. Since there is no established effect measure or outcome for either study, we redo this analysis for survival and solvency. Results: In a case study on bankruptcy following melanoma treatment, the relative risk of bankruptcy suggested that 80-89-year-old patients suffer more financially from melanoma-related expenses than 20-34-year-old patients. The relative risk of solvency and the risk difference suggested the opposite conclusion. Conclusion: To increase transparency around this paradox, researchers reporting one outcome should note if considering the opposite outcome would yield different conclusions. Researchers should also report or estimate underlying risks alongside effect measures when possible.

5. An iterative method for detecting outlying studies in meta-analysis

Zhuo Meng^{1*}, Jingshen Wang², Lifeng Lin³, Chong Wu⁴

¹Department of Statistics, Florida State University

²University of California, Berkeley

³University of Arizona

⁴University of Texas MD Anderson Cancer Center

Meta-analysis is a widely used tool for synthesizing results from multiple studies. The collected studies are deemed heterogeneous when they do not share a common underlying effect size; thus, the factors attributable to the heterogeneity need to be carefully considered. A critical problem in meta-analyses and systematic reviews is that outlying studies are frequently included, which can lead to invalid conclusions and affect the robustness of decision-making. Outliers may be caused by several factors such as study selection criteria, low study quality, small-study effects, etc. Although outlier detection is well-studied in the statistical community, limited attention has been paid to meta-analysis. The conventional outlier detection method in meta-analysis is based on a leave-one-study-out procedure. However, when calculating a potentially outlying study, deviate, other outliers could substantially impact its result. This article proposes an iterative method to detect potential outliers, which reduces such an impact that could confound the detection. Furthermore, we adopt bagging to provide valid inference after removing outliers. Based on simulation studies, the proposed iterative method yields smaller bias and heterogeneity after performing a sensitivity analysis of removing the identified outliers. It also provides higher accuracy on outlier detection. Two case studies are used to illustrate the proposed method in real-world performance.

6. An Association Test for Sequencing Data Based on Kernel Neural Networks

Tingting Hou, Department of Biostatistics, University of Florida

The recent development of artificial intelligence in genomics holds great promise to unravel the complex relationships between genetic variants and disease phenotypes and improve our understanding of the genetic etiology of complex diseases. However, due to the complexity of neural networks and their unknown limiting distributions, building significance tests on neural networks to examine complex genotype-phenotype relationships remains a great challenge. We previously developed a kernel-based neural network (KNN) method, which inherits features from both linear mixed models (LMM) and classical neural networks. Based on the KNN framework, we propose a Wald-type test to evaluate the joint association of a set of genetic variants with a disease phenotype, considering non-linear and non-additive effects. In addition, we also provide two tests to evaluate the linear genetic effect and non-linear/non-additive genetic effects (e.g., interaction effects), respectively. Through simulations, we demonstrated that our proposed method attained higher power compared to the sequence kernel association test (SKAT), especially in the presence of non-linear and interaction effects.

7. Skew-Normal Classification in High-Dimensional data

Haesong Choi^{1*} and Qing Mai¹

¹Department of Statistics, Florida State University

A considerable number of studies have been devoted to high-dimensional classification models under the assumption of normality. However, it may be too restrictive in data analysis. Motivated by the data set that exhibits asymmetry, including environmental, financial, and biomedical ones, we propose a high-dimensional discriminant analysis model called the SKNC model (short for SKew-Normal Classification). By incorporating the skew-normal model, the SKNC model inherits convenient formal properties of the normal distribution and improves its flexibility on skewed data in classification. Theoretical results rigorously show that the SKNC model achieves variable selection, penalized estimation, and prediction consistency, especially in high-dimensional settings. We empirically demonstrate the superior performance of the SKNC model over existing methods in simulated and real datasets.

8. Multimodal Functional Deep Learning for Multi-omics Data

Yuan Zhou^{1*}, Dr. Shan Zhang², Dr. Pei Gang³, Dr. Qing Lu¹

¹Department of Biostatistics, University of Florida

²Michigan State University

³Illinois State University

With rapidly evolving high-throughput technologies and ever-decreasing costs, it becomes feasible to collect diverse types of omics data in large-scale studies. While the multi-omics data generated from these studies hold great promise for innovative insights on biology mechanisms of human disease, the high-dimensionality of omics data and the complexity between various levels of omics data and disease phenotypes bring tremendous analytic challenges. To address these challenges and to facilitate ongoing multi-omics analysis, we propose a Multimodal Functional Deep Learning (MFDL) method for high-dimensional multi-omics data analysis. MFDL model the complex relationships between genetic variants and disease phenotypes through the hierarchical structure of deep neural networks and handle high-dimensional omics data by using the functional data analysis technique. Moreover, MFDL utilizes the structure of the multimodal model to model interactions between multi-omics data. Through simulation studies and real data applications, we demonstrate the advantages of MFNN in terms of prediction accuracy as well as being robust to the high dimensionality and noise of the data.

9. Selection of Variance Component for KNN

Heng Ge^{1*}, Xiaoxi Shen², Qing Lu¹

¹Department of Biostatistics, University of Florida

²Texas State University

Linear mixed model (LMM) and its extensions have long been the state-of-the-art choice for genetic risk prediction analysis of complex diseases. However, LMM commonly assumes a linear genotype-phenotype relationship, which may not be satisfied for diseases involving complex genetic etiology (e.g., interactions). Moreover, it remains computationally challenging to model a large number of samples, especially with multiple kernels built on genetic data. To address these rising challenges, we propose a penalized kernel neural network (KNN) method with a minimum quadratic unbiased estimator (MINQUE). The new method uses the hierarchical structure of kernel neural networks to model complex genotype-phenotype relationships (e.g., non-linear and non-additive relationships), and a closed-form solution is derived to improve computational efficiency. We further extend Closed-form MINQUE to a numerical solution, which can accommodate a variety of penalty terms. Through simulations and analysis PET-imaging outcomes, we have shown that penalized KNN with MINQUE estimator outperformed current state-of-the-art methods in terms of accuracy and statistical efficiency.

10. Inferring cell-cell communications from spatially resolved transcriptomics data

Dongyuan Wu^{1*}, Jeremy Gaskins², Michael Sekulam², Susmita Datta¹

¹Department of Biostatistics, University of Florida

²University of Louisville

Cellular communication through biochemical signaling is fundamental to every biological activity. Investigating cell signaling diffusions across cell types can further help understand biological mechanisms. In recent years, it has become a rather important research topic with the maturity of single-cell sequencing technologies. However, cell signaling activities are spatially constrained, and the single-cell data cannot provide spatial information for each cell. This issue may cause a high false discovery rate, and using spatial transcriptomics data is necessary. On the other hand, as far as we know, most existing methods focus on providing an arbitrary measurement to estimate intercellular communication instead of relying on a statistical model. It is undeniable that descriptive statistics are straightforward and more accessible, while a suitable statistical model can provide a more accurate and reliable estimate. In this way, we propose a generalized linear regression model to infer cellular communications based on spatially resolved transcriptomics data, especially spot-based data. Our Bayesian approach estimates the communication scores between cell types with the consideration of their corresponding distances. Due to the property of the regression model, the proposed method naturally provides the direction of the effect of cell types on a specific signaling pathway that other approaches cannot offer. We conduct simulation studies to assess the performance under different scenarios. We also employ the proposed model in a real-data application and compare it with other existing algorithms. In summary, our innovative model

can fill in gaps in the inference of cell-cell communication and provide a robust and straightforward result.

11. A general framework for online monitoring of dynamic networks

Yipeng Wang^{1*}, Xiulin Xie¹, and Peihua Qiu¹

¹Department of Biostatistics, University of Florida

Dynamic networks where nodes or edges are in constant change become increasingly common due to advances in information technology and data collection systems. A dynamic network process takes the form of a data stream in the sense that new networks keep being collected over time. One fundamental research problem is to monitor the underlying sequential process of the observed networks to see whether it is longitudinally stable, and a signal should be given as soon as an anomaly occurs. To deal with an anomalous change without specializing change types, we propose (a) characterizing each network by correlated summary statistics from the whole network or connected components, and (b) using a nonparametric multivariate cumulative sum (CUSUM) chart to monitor the serially correlated process of the summary statistics. If nodes keep changing in the process, the number of nodes, average degree, the number of connected components, and average diameter are monitored simultaneously. As a special case, if nodes are fixed, the aforementioned summary statistics except the number of nodes are used to be monitored. Traditionally, a monitoring strategy is based on a fixed node set and assumptions for the dynamic network model, which may not be valid in practice. For example, fitting a stochastic block model that requires network communities remain unchanged is questionable when dynamic networks are constantly changing. As a result, monitoring the model-based metrics is dubious. In addition, existing methods based on conventional control charts suffer from unrealistic assumptions of normality or time-independent. The new general framework is robust to use if these assumptions are violated. A procedure based on the Poisson hurdle model is proposed to generate sparse dynamic networks rather than simple random graph models. Numerical studies show that our monitoring scheme is generally more reliable and effective than the representative existing methods.

12. High-dimensional Posterior Consistency in Multi-response Regression models with Non-informative Priors for Error Covariance Matrix

Partha Sarkar^{1*}, Kshitij Khare¹, Malay Ghosh¹

¹Department of Statistics, University of Florida

The Inverse-Wishart (IW) distribution is a standard and popular choice of priors for covariance matrices and has attractive properties such as conditional conjugacy. However, the IW family of priors has crucial drawbacks, including the lack of effective choices for non-informative priors. Several classes of priors for covariance matrices that alleviate these drawbacks, while preserving computational tractability, have been proposed in the literature. These priors can be

obtained through appropriate scale mixtures of IW priors. However, the high-dimensional posterior consistency of models which incorporate such priors has not been investigated. We address this issue for the multi-response regression setting (q responses, n samples) under a wide variety of IW scale mixture priors for the error covariance matrix. Posterior consistency and contraction rates for both the regression coefficient matrix and the error covariance matrix are established in the "large q , large n " setting under mild assumptions on the true data-generating covariance matrix and relevant hyperparameters. In particular, the number of responses q_n is allowed to grow with n , but with $q_n = o(n)$. Also, some results related to the inconsistency of posterior mean for $q_n/n \rightarrow \gamma$, where $\gamma \in (0, \infty)$ are provided.

13. Multivariate single index modeling of longitudinal data with multiple responses

Zibo Tian^{1*} and Peihua Qiu¹

¹Department of Biostatistics, University of Florida

In medical studies, composite indices and/or scores are routinely used for predicting medical conditions of patients. These indices are usually developed from observed data of certain disease risk factors, and it has been demonstrated in the literature that single index models can provide a powerful tool for this purpose. In practice, the observed data of disease risk factors are often longitudinal in the sense that they are collected at multiple time points for individual patients, and there are often multiple aspects of a patient's medical condition that are of our concern. However, most existing single index models are developed for cases with independent data and a single response variable, which are inappropriate for the problem just described in which within-subject observations are usually correlated and there are multiple mutually correlated response variables involved. This paper aims to fill this methodological gap by developing a single index model for analyzing longitudinal data with multiple responses. Both theoretical and numerical justifications show that the proposed new method provides an effective solution to the related research problem. It is also demonstrated using a dataset from the English Longitudinal Study of Aging.

14. Effective Comparison of Two Potentially Crossing Hazard Rate Curves

Xiaoxi Zhang^{1*}, Peihua Qiu¹, Somnath Datta¹

¹Department of Biostatistics, University of Florida

In survival data analysis, comparison of two hazard rate curves is critically important for evaluating a treatment effect. In many applications, the two hazard curves could potentially cross each other, violating the proportional hazards assumption in Cox's model. In such cases, the traditional test like the log-rank test that was developed based on that assumption would be ineffective. There have been discussions in the literature on comparison of two potentially crossing hazard curves, based on either parametric modeling or nonparametric testing approaches. However, the assumed models of the parametric methods are difficult to justify in practice. On the other hand, the nonparametric tests are usually based on maximization with

respect to an unknown crossing point, leading to complex null distributions for the corresponding test statistics. We suggest a novel method for comparing two hazard curves based on a nonparametric testing procedure. Its test statistic avoids the maximization mentioned above and consequently has the desirable asymptotic normality property under some regularity conditions. We show that the new method is effective for comparing two crossing hazard curves.

15. Incorporating Subsampling into Bayesian Model for High-Dimensional Spatial Data

Sudipto Saha^{1*}, Jonathan Bradley¹

¹Department of Statistics, Florida State University

Additive spatial statistical models with weakly stationary process assumptions have become standard in spatial statistics. However, one disadvantage of such models is the computation time, which rapidly increases with the number of datapoints. The goal of this article is to apply an existing subsampling strategy to standard spatial additive models and to derive the spatial statistical properties. We call this new strategy the "spatial data subset model" approach, which can be applied to big datasets in a computationally feasible way. Our approach has the advantage that one does not require any additional restrictive model assumptions (i.e., computational gains increase as model assumptions are removed). This provides a solution to computational bottlenecks that occur when applying methods such as Kriging to "big data". We provide several properties of this new "spatial data subset model" approach in terms of moments, sill, nugget, and range under several sampling designs. We present the results of the "spatial data subset model" approach on a simulated dataset, and on a large dataset consists of 150,000 observations of daytime land surface temperatures measured by the MODIS instrument onboard the Terra satellite.

16. Modeling of clustered censored survival data with spatial correlated random effects using SBART

Durbadal Ghosh^{1*}, Dr Debajyoti Sinha¹, Jonathan Bradley¹, George Rust¹

¹Department of Statistics, Florida State University

This article presents a nonparametric Bayesian survival analysis of clustered survival data where cluster effects are spatially correlated, and some of the spatially-correlated cluster-level covariates are estimated using a data source different from the survival data source. Widespread parametric and semi-parametric hazards regression models for clustered survival data are inappropriate and inadequate for this situation, with complex spatial effects on both clustering and covariates. In this article, we present a general nonparametric model for such clustered censored survival data with spatial correlation under a paradigm of Bayesian ensemble learning called Soft Bayesian Additive Regression Trees or SBART. Our computationally feasible method can incorporate unknown functional forms of the main effects and interactions of various covariates and cluster-specific effects. We illustrate the practical

implementation of our method with an analysis of the effects of various intervenable and non-intervenable covariates on survival times of breast cancer patients from different counties (clusters) in Florida. For our analysis, the clustered survival data with patient-level covariates come from Florida Cancer Registry (FCR), and the data for one county-level intervenable covariates come from the Behavioural Risk Factor Surveillance Survey (BRFS). We compare our method with existing analysis methods to demonstrate our advantage in assessing the impacts of intervention in some cluster/county level and patient-level covariates to eliminate racial disparity in breast-cancer survival in different Florida counties.

Posters:

1. Iterative Multivariate Random Forest for Feature Selection in Integrating Multi-Omics Datasets

Wei Zhang¹ and X. Steven Chen¹

¹Division of Biostatistics, Department of Public Health Sciences, University of Miami, Miller School of Medicine

Integrating multi-omics datasets requires the selection of important features from one data type that are associated with another data type. Traditional methods such as sparse canonical correlation analysis (sCCA) and sparse partial least squares (sPLS) face challenges in handling non-linear, mixed types of information, and interactions present in omics data. Tree-based methods such as random forest are more effective in handling such data. In this work, we propose an iterative multivariate random forest (iterMRF) approach to extract major contributing features between omics datasets. The method uses an iterative framework to select important features in the random forest with a multivariate splitting rule. The feature selection procedures were evaluated by several criteria, including permutation importance measures and variable splitting weights. Our approach shows competitive performance in selecting important features between data types through both simulation studies and real data cases.

2. Future of Data Science Emerging Technology Trends: A perspective analytics on large dataset of research databases

Muhammad Usman Aslam

Department of Science and Engineering (Data Science), University of West Florida

This undertaking analyses the evolution of Emerging Technologies around the globe, using the predictions made in the Horizon Report, published yearly from 2004, MIT Technology review and predictions made by Institute of Electrical and Electronics Engineers (IEEE). This research applies social evaluation, primarily based on Google Trends, and Bibliometric analysis, with data

of scientific publications from IEEE, MIT, Questia online Library, JSTOR, Springer, Hindawi and Web of Science, with a purpose to discover which technology had been a hit and sincerely impacted mainstream training, and which one failed to have the anticipated impact. This mission gives guidelines that can be beneficial to those who are looking forward to investing in new research regions.

3. Data-driven Evaluation of Trajectories in Single-Cell RNA-Seq Data

Xiaoru Dong¹ and Rhonda Bacher¹

¹Department of Biostatistics, University of Florida

Trajectory inference (TI) is an important component in single cell RNA-sequencing data, which orders single cells along trajectories based on the similarity of expression and consequently to model the process of dynamic cellular changes. However, a great number of preliminary tasks such as feature selection and dimension reduction need to be completed before proceeding to TI. We found that the choice of parameters in these preprocessing steps can greatly affect the reliability of the trajectory estimation. We created a novel 2-phase framework to evaluate what parameters should be used in preprocessing stage with the aim of obtaining more accurate trajectories in the TI analysis. Our 2-phase framework consists of detecting the existence of trajectories and evaluating the goodness of trajectory fitting, which assess credibility of the estimated trajectory in the high and low dimensional spaces based on the similarity between cells respectively.

4. Novel Spatial Two-Fold Modeling of Perimetric Rates of Change Using Glaucomatous Population Data

Chen Zhao^{1*}, J. Sunil Rao¹, Swarup S. Swaminathan²

¹Public Health Science, University of Miami

²University of Miami Health System, Bascom Palmer Eye Institute

Static automated perimetry is the most common form of measuring visual fields. In glaucomatous eyes, it is generally assumed that visual loss can be due to either a global loss or focal loss of visual sensitivities. In this work, we will consider modeling focal visual loss over time in a series of 500 glaucoma patients at the Bascom Palmer Eye Institute. Specifically, we have repeated visual field tests for each individual along with information on three different types of focal clustering patterns. The goal is to better identify fast progressors amongst this group of patients. To model this data, we develop two-fold random slope models which can be seen as an extension of the two-fold small area estimation models of Torabi and Rao (2014). We also include spatial structure in the fixed effects to account for the orientation of focal clusters. Our model borrows strength across individuals and across connected focal clusters resulting in more accurate estimates (as measured by estimated MSPE) of individual rates of visual sensitivity decline. Confidence intervals are estimated using parametric bootstrapping.

5. Neural-network transformation models for counting processes

Rongzi Liu¹, Chenxi Li², Qing Lu¹

¹Department of Biostatistics, University of Florida

²Michigan State University

While many survival models have been invented, the Cox model and the proportional odds model are among the most popular ones. Both models are special cases of the linear transformation model. The linear transformation model typically assumes a linear function on covariates, which may not reflect the complex relationship between covariates and survival outcomes. Nonlinear functional form can also be specified in the linear transformation model. Nonetheless, the underlying functional form is unknown and mis-specifying it leads to biased estimates and reduced prediction accuracy of the model. To address this issue, we develop a neural-network transformation model. Similar to neural networks, the neural-network transformation model uses its hierarchical structure to learn complex features from simpler ones and is capable of approximating the underlying functional form of covariates. It also inherits advantages from the linear transformation model, making it applicable to both time-to-event analyses and recurrent event analyses. Simulations demonstrate that the neural-network transformation model outperforms the linear transformation model in terms of estimation and prediction accuracy when the covariate effects are nonlinear. The advantage of the new model over the linear transformation model is also illustrated via two real applications.

6. Assessing the impact of key parameters on the results of meta-analyses

Wenshan Han, Department of Statistics, Florida State University

Systematic reviews and meta-analyses are essential tools in contemporary evidence-based medicine. They have been extensively used to synthesize diverse sources of evidence and draw a conclusive and unified decision for medical decision-making. However, in recent years, conflicting conclusions from different meta-analyses conducted by separate teams on the same research topic have raised concerns about the reliability of these studies. One reason for this issue is the sensitivity of meta-analysis results to certain factors, such as the inclusion criteria for eligible studies in the systematic review and the choice of meta-analysis models. This essay aims to evaluate the conclusion of a meta-analysis using the arm-based model under the Bayesian framework. The arm-based meta-analysis model focuses on estimating the marginal results of treatment arms and provides flexibility for drawing conclusions with different effect measures. Despite its benefits, the arm-based model heavily relies on the correlation between treatment groups, a crucial parameter that accounts for randomization in clinical trials. Moreover, the between-study variance parameters play a critical role in meta-estimates by reflecting heterogeneity. Our proposed methods employ the spirit of tipping point analysis, a widely adopted technique in missing data imputation, to evaluate the robustness of meta-

analysis results concerning the estimation of key parameters. Additionally, we introduce innovative visualization tools that intuitively display the impact of key parameters on meta-results. We apply these methods to two real-world meta-analyses and observe gradual changes in meta-analysis results as we assign key parameters to a series of values.

7. Single Cell Linear Adaptive Negative-binomial Expression Testing (scLANE)

Jack R. Leary¹ and Rhonda Bacher¹

¹Department of Biostatistics, University of Florida

Single cell RNA-sequencing (scRNA-seq) offers a high-resolution view of cellular biology, including dynamic processes such as differentiation and disease progression. Many methods have emerged that estimate a cell-level time ordering from static scRNA-seq samples, which use similarity of gene expression to place cells in order on some biological manifold. Researchers then typically 1) assume the ordering represents a biological process and 2) characterize which genes are associated with that process. Changes in expression over trajectories are usually complex with generalized additive models (GAMs) the dominant current choice of model. However, while GAMs are excellent for fitting nonlinear relationships, they are not easily interpretable. To address this trade-off, we developed Single Cell Linear Adaptive Negative-binomial Expression (scLANE) testing. scLANE balances the need for a flexible nonlinear model and the facilitation of biological interpretation. We demonstrate our method's accuracy and ability to draw meaningful comparisons from the model on simulated data and a case-study datasets having tens of thousands of cells and from multiple subjects.

8. Cherry Blossom Prediction

Nitul Singha, Department of Mathematics & Statistics, University Of West Florida

The blooming of cherry blossoms is a highly anticipated event in many parts of the world, attracting tourists and locals alike to witness the beauty of these delicate flowers. However, predicting the arrival of cherry blossoms has become a complex and fascinating challenge, as weather patterns and climate change can greatly impact the timing of their blooms. In 2023, the American Statistical Association (ASA) is hosting a cherry blossom prediction competition, inviting students, researchers, and citizen scientists to predict the peak bloom date of cherry trees in four locations around the world: Washington D.C., USA; Kyoto, Japan; Vancouver, Canada; and Liestal-Weideli, Switzerland.

The competition presents a unique opportunity for participants to apply their skills in predictive modeling to a real-world problem, as predicting cherry blossom blooms has important implications for tourism, agriculture, and ecology. The data available for the competition dates back to 1981, allowing participants to explore patterns in cherry blossom phenology over the past few decades. However, accurately predicting the timing of cherry blossom blooms remains

a significant challenge due to the complexity of weather patterns and the impact of climate change on plant phenology.

To address this challenge, we will be using two common predictive modeling techniques in our project: logistic regression and linear regression. Logistic regression is a widely used statistical technique that is particularly effective for binary outcomes, such as the "month" of cherry blossoms. Linear regression, on the other hand, is a powerful tool for predicting numerical outcomes, such as the timing of cherry blossom blooms. By analyzing historical data from the four target cities and applying these two techniques, we hope to develop accurate models that can make reliable predictions for the 2023 competition.

Our project will contribute to the scientific understanding of cherry blossom phenology and the impact of climate change on plant growth and development. The findings of our study may also have practical applications in fields such as agriculture and tourism, where accurate predictions of cherry blossom blooms can have important economic and cultural implications. We believe that our project will help further the development of predictive modeling techniques and provide valuable insights for future prediction competitions.

9. Customizable SAS Macro Program for Case-Control Matching

Kyle Grealis, Department of Public Health Sciences, University of Miami

Case-control matching is a useful technique to understand how a condition or disease affects individual study participants. Each participant with the disease of interest (case) is matched with a preset number of participants without the disease of interest (control) to reduce the effects of confounding. Cases and controls can be matched based on certain characteristics, such as age, gender, socioeconomic status, or numerous comorbid factors. The matching process is susceptible to selection bias and, depending on the size of the dataset or number matching factors to be considered, it can be time intensive. Online statistical software communities offer peer-to-peer support in solving many programming problems but have nonetheless resulted in a patchwork of solutions. Oftentimes the solutions are not translatable to new projects or are invariably difficult to modify and reuse on future projects. The author has improved upon existing case-control sorting and matching programs by providing the end-user with a flexible SAS macro program that reduces selection biases and time required to complete case-control matching. The macro requires five (5) arguments for execution; it allows customization based on age, age range, up to two (2) other matching factors, and the ability to set the case-to-control ratio. The macro matches cases to eligible controls by completing a sorting and randomized matching procedure. Controls are matched to cases only once. The total number of case-control matches achieved is counted and the results are saved. This iterative process is completed 100 times and selects the permutation with the greatest number of matches, i.e., the least number of insufficient matches. Finally, the macro generates two full-report PDF files and saves them to the project folder. One file contains a by-case result with respective controls, while the other lists cases with insufficient matching. This program

streamlines the process from data collection to analysis and simultaneously limits the risk of selection bias. Because it is a macro program, it is reusable across projects without the need to modify and rewrite code. At bottom, users may allocate more time to the analytic processes and investigators can improve study efficiency.