

Preliminary Program for Quantile Regression and Data Heterogeneity Workshop
(Feb 12-13, 2022)

Saturday, February 12

8:15 am Food & Drinks

8:50 am Welcome

9:00-10:20 am Session #1

Regina Liu (Rutgers University): Fusion learning: combine inferences from diverse data sources with heterogeneous data

Ivan Fernandez-Val (Boston University): Dynamic Heterogeneous Distribution Regression Panel Models with an Application to Labor Income Processes

Chair/Discussant: **Roger Koenker**

10:20 am Coffee break

10:40 am - 12:00 pm Session #2

Kean Ming Tan (University of Michigan): Convolution-Type Smoothing Approach for Quantile Regression

Jingshen Wang (UC Berkeley): Debiased Inference on Heterogeneous Quantile Treatment Effects with Regression Rank-Scores

Chair/Discussant: **Lan Wang, Ganggang Xu**

Noon-1:30 pm lunch break

1:30 - 2:50 pm Session #3

Linglong Kong (University of Alberta): Statistical inference for smoothed quantile regression with streaming data

Ying Wei (Columbia University): Integrated Quantile Rank Test (IQRAT) for Gene-level associations

Chair/Discussants: **Emma Zhang, Hongyu Zhao**

2:50-3:10 pm Coffee break

3:10-4:30 pm Session #4 Invited Presentations from Junior Scholars

Shuangning Li (Stanford University): Random Graph Asymptotics for Treatment Effect Estimation under Network Interference

Ben Sherwood (University of Kansas): Computationally efficient penalized quantile regression

Bo Wei (University of Michigan): Estimation of Complier Super-quantile Causal Effects with a Binary Instrumental Variable

Harold Chiang (University of Wisconsin-Madison): Inference for high-dimensional exchangeable arrays

Chair/discussant: **Snigdha Panigrahi**

4:30-4:40 pm break

4:40-5:10 pm FRG group report and discussion

6:00 - 8:00 pm Workshop dinner

Sunday, February 13

9:00 am Food & Drinks

9:30-10:50 am Session #5

Matias Cattaneo (Princeton University): On Binscatter

Bodhi Sen (Columbia University): Measuring association on topological spaces using kernels and geometric graphs

Chair/Discussant: **Kengo Kato**

10:50-11:10 am coffee break

11:10am -12:30pm Session #6

Sunil Rao (University of Miami): A Tour of Classified Mixed Model Predictions and Projections

Xiaohong Chen (Yale University): Adaptive Estimation and Uniform Confidence Bands for Nonparametric IV

Chair/Discussants: **Qi Zheng, Xuming He**

Afternoon: free discussion and social time

Talk Titles and Abstracts

1. Kean Ming Tan (University of Miami, keanming@umich.edu)

Title: Convolution-Type Smoothing Approach for Quantile Regression

Abstract: Quantile regression is a powerful tool for learning the relationship between a response variable and a multivariate predictor while exploring heterogeneous effects. However, the non-smooth piecewise linear loss function introduces challenges to the computational aspect when the number of covariates is large. To address the aforementioned challenge, we propose a convolution-type smoothing approach that turns the non-differentiable quantile piecewise linear loss function into a twice differentiable, globally convex, and locally strongly convex surrogate, which admits a fast and scalable gradient-based algorithm to perform optimization. In the low-dimensional setting, we establish nonasymptotic error bounds for the resulting smoothed estimator. In the high-dimensional setting, we propose the concave regularized smoothed quantile regression estimator, which we solve using a multi-stage convex relaxation algorithm. Theoretically, we characterize both the algorithmic error due to non-convexity and statistical error for the resulting estimator simultaneously. We show that running the multi-stage algorithm for a few iterations will yield an estimator that achieves the oracle property. Our results suggest that the smoothing approach leads to a significant computational gain without a loss in statistical accuracy.

2. Ivan Fernandez-Val (Boston University, ivanf@bu.edu)

Title: "Dynamic Heterogeneous Distribution Regression Panel Models with an Application to Labor Income Processes," joint with Wayne Gao, Yuan Liao and Francis Vella

Abstract: We consider estimation of a dynamic distribution regression panel data model with heterogeneous coefficients across units. The objects of interest are functionals of these coefficients including linear projections on unit level covariates. We also consider actual, stationary and counterfactual distributions of the outcome variable. We investigate how changes in initial conditions or covariate values affect these objects. Coefficients and their functionals are estimated via fixed effect methods, which are debiased to deal with the incidental parameter problem. We propose a cross-sectional bootstrap scheme to perform uniform inference on

function-valued objects. This avoids coefficient re-estimation and is shown to be consistent for a large class of data generating processes, including the reference point of coefficient homogeneity (conditional on covariates). We employ annual labor income data from the PSID to illustrate the variety of empirical issues we can address. First, we predict the impact of hypothetical tax policies and find substantially smaller effects than those from models based on homogeneous autoregressive and distributional regression processes. Second, we examine the impact on the distribution of labor income from increasing the education level of a chosen subsample of workers. Explicitly, we increase the education level of all workers with less than 12 years of schooling to that level of schooling and find short and long run increases in the distribution at the bottom tail with the upper tail relatively unaffected. Finally we uncover notable heterogeneity in income mobility implying substantial individual heterogeneity in the incidence to be trapped in poverty.

3. Jingshen Wang (University of California at Berkeley, jingshenwang@berkeley.edu)

Title: Debiased Inference on Heterogeneous Quantile Treatment Effects with Regression Rank-Scores

Abstract: Understanding treatment effect heterogeneity in observational studies is of great practical importance to many scientific fields because the same treatment may affect different individuals differently. Quantile regression provides a natural framework for modelling such heterogeneity. In this paper, we propose a new method for inference on heterogeneous quantile treatment effects that incorporates high-dimensional covariates. Our estimator combines a L1-penalized regression adjustment with a quantile-specific bias correction scheme based on quantile regression rank scores. We present a comprehensive study of the theoretical properties of this estimator, including weak convergence of the heterogeneous quantile treatment effect process to the sum of two independent, centered Gaussian processes. We illustrate the finite-sample performance of our approach through Monte Carlo experiments and an empirical example, dealing with the differential effect of statin usage for lowering LDL cholesterol levels for the Alzheimer's disease patients who participated in the UK Biobank study. This is joint work with Alexander Giessing at the University of Washington.

4. Bodhi Sen (Columbia University, bodhi@stat.columbia.edu)

Title: Measuring association on topological spaces using kernels and geometric graphs

Abstract: We propose and study a class of simple, nonparametric, yet interpretable measures of association between two random variables X and Y taking values in general topological spaces.

These nonparametric measures -- defined using the theory of reproducing kernel Hilbert spaces - capture the strength of dependence between X and Y and have the property that they are 0 if and only if the variables are independent and 1 if and only if one variable is a measurable function of the other. Further, these population measures can be consistently estimated using the general framework of graph functionals which include k -nearest neighbor graphs and minimum spanning trees. Moreover, a sub-class of these estimators are also shown to adapt to the intrinsic dimensionality of the underlying distribution. Some of these empirical measures can also be computed in near-linear time. Under the hypothesis of independence between X and Y , these empirical measures (properly normalized) have a standard normal limiting distribution. Thus, these measures can also be readily used to test the hypothesis of mutual independence between X and Y . In fact, as far as we are aware, these are the only procedures that possess all the above mentioned desirable properties.

6. Linglong Kong (University of Alberta, lkong@ualberta.ca)

Title: Statistical inference for smoothed quantile regression with streaming data

Abstract: In this paper, we address the problem of how to conduct valid statistical inference for quantile regression with streaming data. The main challenges are that the quantile regression loss function is non-smooth and it is often infeasible to store the entire dataset in memory. To overcome these issues, we propose a fully online updating method for statistical inference in smoothed quantile regression with streaming data. Our main contributions are twofold. First, in the low-dimensional regime, we present an incremental updating algorithm to obtain the smoothed quantile regression estimator with streaming data, which allows to construct asymptotically exact statistical inference procedures. Second, in high-dimension regime, we develop an online debiased LASSO procedure to accommodate the special sparse structure of streaming data. Our procedure is updated with only the current data and summary statistics of historical data, and corrects an approximation error from online updating with streaming data. Moreover, theoretical results such as estimation consistency and asymptotic normality are established to justify its validity in both settings. Numerical studies including extensive simulations and a real data example confirm that the proposed method performs well in practical settings.

7. Ying Wei (Columbia University, yw2148@cumc.columbia.edu)

Title: Integrated Quantile Rank Test (IQRAT) for Gene-level associations

Abstract: Gene-based testing is a commonly employed strategy in genetic association studies. Gene-trait associations are complex due to underlying population heterogeneity, gene-environment interactions, and various other reasons. Existing gene-based tests, such as Burden and Sequence Kernel Association Tests (SKAT), focus on mean-level associations and may miss or underestimate higher-order associations that could be scientifically interesting. We introduce a new family of gene-level association tests that integrate the quantile-rank score process to accommodate complex associations better. The resulting test statistics enjoy multiple advantages. First, they are almost as efficient as the best existing tests when the associations are homogeneous across quantile levels and have improved efficiency for complex and heterogeneous associations. Second, they provide valuable insights into risk stratification. Third, the test statistics are distribution-free and could accommodate a wide range of underlying distributions; We established the asymptotic properties under the null and alternative hypotheses to validate the proposed tests theoretically. We also conducted extensive simulations to assess its empirical performance compared to the existing approaches. Finally, we illustrate its real-world applications to identify genes associated with lipid traits using a Metabochip dataset and identifying eGenes from the multi-tissue gene-expression data in GTEx.

8. J. Sunil Rao (University of Miami, JRao@med.miami.edu)

Title: A Tour of Classified Mixed Model Predictions and Projections

Abstract: Many practical problems are related to prediction where the main interest is at the subject or small sub-population level. In such cases, it's possible to make substantial gains in prediction accuracy by identifying a class that a new subject belongs to and associating the new subject with a random effect corresponding to the same class in the training data so that mixed model prediction can be used. In this talk, we first introduce the original classified mixed model prediction idea and then discuss some newer developments in multivariate classified mixed model prediction and classified mixed model projections for new data outside the range of the training data. This is joint work over the years with Jiming Jiang (UC-Davis), Thuan Nguyen (OHSU) and former UM students Menyng Li (Moderna), Hang Zhang (Biogen) and Jie Fan (Novartis).

11. Shuangning Li (Stanford University, lsn@stanford.edu)

Title: Random Graph Asymptotics for Treatment Effect Estimation under Network Interference

Abstract: The network interference model for causal inference places all experimental units at

the vertices of an undirected exposure graph, such that treatment assigned to one unit may affect the outcome of another unit if and only if these two units are connected by an edge. This model has recently gained popularity as means of incorporating interference effects into the Neyman--Rubin potential outcomes framework; and several authors have considered estimation of various causal targets, including the direct and indirect effects of treatment. In this paper, we consider large-sample asymptotics for treatment effect estimation under network interference in a setting where the exposure graph is a random draw from a graphon. When targeting the direct effect, we show that---in our setting---popular estimators are considerably more accurate than existing results suggest, and provide a central limit theorem in terms of moments of the graphon. Meanwhile, when targeting the indirect effect, we leverage our generative assumptions to propose a consistent estimator in a setting where no other consistent estimators are currently available. We also show how our results can be used to conduct a practical assessment of the sensitivity of randomized study inference to potential interference effects. Overall, our results highlight the promise of random graph asymptotics in understanding the practicality and limits of causal inference under network interference. This is joint work with Stefan Wager.

12. Ben Sherwood (ben.sherwood@ku.edu, University of Kansas)

Title: Computationally efficient penalized quantile regression

Abstract: Quantile regression with a lasso penalty can be framed as a linear programming problem. If a group lasso penalty is used, then it becomes a second order cone programming problem. These approaches become computationally burdensome for large values of n or p . Using a Huber approximation to the quantile function allows for the use of computationally efficient algorithms that require a differentiable loss function that can be implemented for both penalties. These algorithms then can be used as the backbones for implanting penalized quantile regression with other penalties such as Adaptive Lasso, SCAD, MCP and group versions of these penalties.

13. Bo Wei (boweinju@umich.edu, University of Michigan)

Title: Estimation of Complier Super-quantile Causal Effects with a Binary Instrumental Variable

Abstract: Estimating causal effect of a treatment or exposure for an interested subpopulation is a fundamental interest in many biomedical and economical studies. Super-quantile, also known as expected shortfall, is an attractive measure of risk for the subpopulation because it can capture heterogeneity and aggregated local information of effect over a range of distribution of outcomes

simultaneously. In this work, we propose a complier super-quantile causal effect (CSQCE) model under instrumental variable (IV) framework to quantify the CSQCE for the data with unmeasured confounders. By utilizing the special characteristic of binary IV, we propose a simple and easily-implemented two-step estimation procedure, which can simply be solved by weighted linear regression and weighted quantile regression. We rigorously justify the asymptotic properties for the proposed estimator. Extensive simulations have been conducted to confirm its validity and satisfactory finite-sample performance. An application to a dataset from National Job Training Partnership Act (JTPA) study demonstrates the practical utility of the proposed method.

14. Harold Chiang (hdchiang@wisc.edu, University of Wisconsin-Madison.)

Title: Inference for high-dimensional exchangeable arrays

Abstract: We consider inference for high-dimensional separately and jointly exchangeable arrays where the dimensions may be much larger than the sample sizes. For both exchangeable arrays, we first derive high-dimensional central limit theorems over the rectangles and subsequently develop novel multiplier bootstraps with theoretical guarantees. These theoretical results rely on new technical tools such as Hoeffding-type decomposition and maximal inequalities for the degenerate components in the Hoeffding-type decomposition for the exchangeable arrays. We exhibit applications of our methods to uniform confidence bands for density estimation under joint exchangeability and penalty choice for ℓ_1 -penalized regression under separate exchangeability. Extensive simulations demonstrate precise uniform coverage rates. We illustrate by constructing uniform confidence bands for international trade network densities.