

*2015 Joint Statistical Workshop of
The Chinese University of Hong Kong
and Academia Sinica*

April 19th- 20th, 2015

Department of Statistics, The Chinese University of Hong Kong

Institute of Statistical Science, Academia Sinica

2015 JOINT STATISTICAL WORKSHOP

of The Chinese University of Hong Kong and Academia Sinica

April. 19 (Sunday)

09:00-09:30	Opening: Ching-Shui Cheng 鄭清水 ● Qi Man Shao 邵啟滿 ● Hoi Ying Wong 王海嬰 ●
	Session 1 Chair: Phillip Yam 任尚智 ●
09:30-10:10	Tony Sit 薛賢鴻 ● Accelerated failure time model under general biased sampling scheme
10:10-10:50	Hsin-wen Chang 張馨文 ● Tests for stochastic ordering under biased sampling
10:50-11:20	Coffee Break
	Session 2 Chair: Ci-Ren Jiang 江其衽 ●
11:20-12:00	Yuan Yuan Lin 林媛媛 ● Nearly semiparametric efficient estimation of quantile regression
12:00-12:25	Yifan Wang 王亦璠 ● Bayesian quantile structural equation models
12:25-14:00	Lunch Break
	Session 3 Chair: Ying Ying Wei 魏穎穎 ●
14:00-14:40	Hsin-Chou Yang 楊欣洲 ● Recovering the missing heritability of complex diseases
14:40-15:20	Xiao-Dan Fan 樊曉丹 ● Model-based inference of protein dynamics based on MD data
15:20-15:50	Coffee Break
	Session 4 Chair: Hsin-wen Chang 張馨文 ●
15:50-16:30	Ci-Ren Jiang 江其衽 ● Inverse regression for longitudinal data
16:30-16:55	Yang An 安楊 ● Short-term stock price prediction based on limit order book dynamics
16:55-17:20	Ying Wang 王瑩 ● Bayesian option pricing framework with stochastic volatility
17:20-17:45	Yuan-Lung Lin 林遠隆 ● General supplementary difference sets (GSDS): A key to the construction of (near)-Hadamard designs
18:00	Bus leaves Institute of Statistical Science at 18:00
18:30-20:30	Reception (By invitation only) Hao Liao Li Chinese Restaurant 好料理麗緻喜宴

April. 20 (Monday)

	Session 5 Chair: Xiao-Dan Fan 樊曉丹 ●
09:00-09:30	Tai-Chi Wang 王泰期 ● Generalized framework for detecting communities of social networks by the scanning method
09:30-10:10	Chun Yip Yau 邱俊業 ● Inference for multiple change-points in time series via likelihood ratio scan statistics
10:10-10:50	Frederick Kin Hing Phoa 潘建興 ● The swarm intelligence based (SIB) method and its application in statistics
10:50-11:20	Coffee Break
	Session 6 Chair: Tso-Jung Yen 顏佐榕 ●
11:20-12:00	Ying Ying Wei 魏穎穎 ● Correlation motif model for integrative analyses of genomic data
12:00-12:25	Ping Yang 楊平 ● Multiple comparisons with two controls for ordered categorical responses
12:25-14:00	Lunch Break
	Session 7 Chair: Frederick Kin Hing Phoa 潘建興 ●
14:00-14:40	Ting-Li Chen 陳定立 ● On the weak convergence and central limit theorem of blurring and nonblurring processes with application to robust location estimation
14:40-15:20	Phillip Yam 任尚智 ● Globally efficient nonparametric inference of average treatment effects by empirical balancing calibration weighting
15:20-15:50	Coffee Break
	Session 8 Chair: Yuan Yuan Lin 林媛媛 ●
15:50-16:30	Tso-Jung Yen 顏佐榕 ● Parameter clustering: encouraging similarities between estimates via euclidean distance regularization
16:30-16:55	Zheng Zhang 張政 ● Well-posedness of mean-field type forward-backward stochastic differential equations
16:55-17:20	Jie Hu 胡杰 ● Bayesian detection of embryonic gene expression onset in C. Elegans
17:20-17:45	Shih-Kai Chu 朱是鏞 ● Genome-wide pharmacogenomic study on methadone maintenance treatment identifies SNP rs17180299 and multiple haplotypes on <i>CYP2B6</i> , <i>SPON1</i> , and <i>GSG1L</i> associated with plasma concentrations of methadone R- and S-enantiomers in heroin-dependent patients
18:00	Bus leaves Institute of Statistical Science at 18:00
18:30-20:30	Banquet (By invitation only) Kuang Wu Restaurant 光武農產品生活館

Accelerated failure time model under general biased sampling scheme

Tony Sit

Department of Statistics, The Chinese University of Hong Kong

Abstract

Right-censored time-to-event data are sometimes observed from a (sub) cohort of patients whose survival times can be subject to outcome-dependent sampling schemes. In this paper, we propose a unified estimation method for semiparametric accelerated failure time models under general biased estimating schemes. The proposed estimator of the regression covariates is developed upon a bias-offsetting weighting scheme and is proved to be consistent and asymptotically normally distributed. Large sample properties for the estimator are also derived. Using rank-based monotone estimating functions for the regression parameters, we find that the estimating equations can be easily solved via convex optimisation. The methods are confirmed through simulations and illustrated by application to real data sets on various sampling scheme including length-bias sampling, the case-cohort design and its variants. This is a joint work with Jane Paik Kim and Zhiliang Ying.

Tests for stochastic ordering under biased sampling

Hsin-wen Chang

Institute of Statistical Science, Academia Sinica

Abstract

In two-sample comparison problems it is often of interest to examine whether one distribution function majorizes the other, i.e. for the presence of stochastic ordering. This talk introduces a nonparametric test for stochastic ordering based on size-biased data, allowing the pattern of size bias to differ between the two samples. The test is formulated in terms of a maximally-selected local empirical likelihood statistic. A Gaussian multiplier bootstrap is devised to calibrate the test. A simulation study indicates that the proposed test outperforms an analogous Wald-type test, and that it provides substantially greater power than what is available when ignoring the sampling bias. The approach is illustrated using data on blood alcohol concentration and age of drivers involved in car accidents, in which size bias is present because the drunker drivers are more likely to be sampled. Further, younger drivers tend to be more affected by alcohol, so when comparing with older drivers, the analysis is adjusted for differences in the patterns of size bias.

Nearly semiparametric efficient estimation of quantile regression

Yuan Yuan Lin

Department of Statistics, The Chinese University of Hong Kong

Abstract

As a competitive alternative to the least squares regression, the quantile regression is a popular statistical tool for the modeling and inference of conditional quantile function. In conventional quantile regression models, major complications involve in the semiparametric efficient estimation arise from density estimation and computational difficulty. This paper proposes a semiparametric efficient estimation and inference procedure for the quantile regression models with global linearity assumption. The basis of the procedure is an efficient score function derived from a least favorable submodel method. The key ingredients are curve estimation and a numerical algorithm. The resulting estimator is proved to be (nearly) semiparametric efficient, with asymptotic variance achieves the semiparametric efficiency lower bound. Numerical studies with supportive evidence are presented.

Bayesian quantile structural equation models

Yifan Wang

Department of Statistics, The Chinese University of Hong Kong

Abstract

Structural equation modeling is a common multivariate technique for the assessment of the interrelationships among latent variables. In recent years, structural equation models (SEMs) have been extensively applied to behavioral, medical, and social sciences. Basic SEMs consist of a measurement equation for characterizing latent variables through multiple observed variables and a mean regression-type structural equation for investigating how explanatory latent variables influence outcomes of interest. However, the mean regression-type structural equation does not provide a comprehensive analysis of the relationship between latent variables. In this paper, we introduce the quantile regression method into SEMs to assess the conditional quantile of the outcome latent variable given the explanatory latent variables and covariates. The estimation is conducted in a Bayesian framework with MCMC algorithm, and the posterior inference is performed with the help of asymmetric Laplace distribution. A simulation study shows that the proposed method performs satisfactorily. An application to a real study on the risk factors of chronic kidney disease is presented.

Recovering the missing heritability of complex diseases

Hsin-Chou Yang

Institute of Statistical Science, Academia Sinica

Abstract

Since 2005, genome-wide association studies (GWAS) have identified single nucleotide polymorphisms (SNPs) associated with complex diseases and traits. According to the NHGRI GWAS Catalog reports, up to the end of 2013, about 12 thousands of disease- or trait-associated SNPs from 1,750 curated publications have been reported. Nevertheless, “missing heritability” of complex diseases and traits is still a gap to be filled. In order to recover the missing heritability in GWAS, we developed statistical and bioinformatics tools to further enrich the collection of disease genes and mutations underlying complex diseases and traits. The developed analysis tools have been applied to study genetic mechanism and disease prognosis of hypertension, rheumatoid arthritis, schizophrenia, and acute lymphoblastic leukemia. The analyses include but are not limited to whole-genome sequencing data analysis, homozygosity disequilibrium analysis, gene- and pathway-based association and interaction analysis, copy number analysis, expression quantitative trait locus mapping, metabolomics data analysis, and integrative omics analysis. The analysis tools and some of their practical applications in biomedical research will be introduced in this talk.

April 19th, 2015
14:40-15:20

Model-based inference of protein dynamics based on MD data

Xiao-Dan Fan

Department of Statistics, The Chinese University of Hong Kong

Abstract

Molecular Dynamics (MD) data is one major source of protein dynamics information, but it is still unclear how to use these data to elucidate the energy landscape of the protein structure space. We proposed a method to integrate thousands of MD trajectories to analyze the protein dynamics. We will use both geometric distance and dynamic distance to probe the grouping structure of conformations within the landscape. A parametric Bayesian approach is used to partition the structural space.

Inverse regression for longitudinal data

Ci-Ren Jiang

Institute of Statistical Science, Academia Sinica

Abstract

Sliced inverse regression is an appealing dimension reduction method for regression models with multivariate covariates. It has been extended to functional covariates where the whole trajectories of random functional covariates are observed completely. We aim at developing an inverse regression approach for intermittently and sparsely measured longitudinal covariates. We show, under some regularity conditions, that the convergence rate of the estimated EDR, effective dimension reduction, directions is a function of smoothing parameters, sample size and the number of eigenfunctions used to reconstruct the empirical inverse of the covariance operator. Simulation studies and data analysis are also provided to demonstrate the performance of our method.

Short-term stock price prediction based on limit order book dynamics

Yang An, Ngai Hang Chan

Department of Statistics, The Chinese University of Hong Kong

Abstract

Interaction of capital market participants is a complicated dynamic process. A stochastic model is proposed to describe the dynamics to predict short-term stock price behaviors. Independent compound Poisson processes are introduced to describe the occurrences of market orders, limit orders, and cancellations of limit orders, respectively. Based on high-frequency observations of the limit order book, the maximum empirical likelihood estimator (MELE) is applied to estimate the parameters of the compound Poisson processes. Moreover, an analytical formula is derived to compute the probability distribution of the first-passage time of a compound Poisson process. Based on this formula, the conditional probability of price increase and the conditional distribution of the duration until the first change in mid-price are obtained. Finally, a novel approach of short-term stock price prediction is proposed and this methodology works reasonably well in the data analysis of Intel (INTC).

Bayesian option pricing framework with stochastic volatility

Ying Wang

Department of Statistics, The Chinese University of Hong Kong

Abstract

The application of stochastic volatility (SV) models in the option pricing literature usually assumes that the market has sufficient option data to calibrate the risk-neutral parameters of the model. When option data are not available, market practitioners have to estimate the model from the historical returns of the underlying asset and then transform the resulting model to its risk-neutral equivalent. However, the likelihood function of an SV model can only be expressed in a high dimensional integration, making the estimation a highly challenging task. Bayesian approach has been the classical way to estimate SV models under the data-generating (physical) probability measure but it is still unclear about the transformation from the estimated physical dynamic to its risk-neutral counterpart. Inspired by the Duan (1995) GARCH option pricing approach, we propose an SV model that enables us to perform Bayesian inference and transformation to risk-neutral dynamic simultaneously and conveniently. Our model relaxes the normality assumption on innovations of both return and volatility processes. Our empirical study shows that the estimated option prices generate a realistic implied volatility smile shapes. In addition, the volatility premium is almost flat across strike prices so that adding a few option data to the historical time series of the underlying asset can greatly improve the estimation of option prices.

General supplementary difference sets (GSDS): A key to the construction of (near)-Hadamard designs.

Yuan-Lung Lin^{*}, Frederick Kin Hing Phoa

Institute of Statistical Science, Academia Sinica

Abstract

We propose a new and unified construction method, namely general supplementary difference sets (GSDS), for near-Hadamard designs when the run sizes are $n \equiv 2 \pmod{4}$. These designs possess high D -efficiencies. Ehlich (1964) derived an upper bound for the determinant of matrices of order $n \equiv 2 \pmod{4}$ achievable only if $2n - 2$ is a sum of two squares. Even in a small range from 1 to 100, there are 6 parameters, 22; 34; 58; 70; 78 and 94, that do not fulfill this condition. We construct these designs for many values of n and formulate a new class of near-Hadamard designs whose determinants are very close to Ehlich's upper bound.

Generalized framework for detecting communities of social networks by the scanning method

Tai-Chi Wang^{*}, Frederick Kin Hing Phoa

Institute of Statistical Science, Academia Sinica

Abstract

With the growth in the big data regime and the popularity of social media, recognizing and analyzing social network patterns are important issues. In real world, society offers a wide variety of possible communities, such as schools, families, and firms. For this reason, community/cluster detection draws much attention as it is important to many applications in business and social sciences. The scan statistics have been verified as a useful tool to determine both structure and attribute clusters in networks. However, most of previous methods assumed that the baseline network model follows the Poisson distribution assumption. In this paper, we generalize the previous scan statistic to accommodate to random connection probability model and logit model. Simulation studies show that the generalized methods have better detection results, and empirical studies show the differences among the proposed methods and the previous methods.

Inference for multiple change-points in time series via likelihood ratio scan statistics

Chun Yip Yau

Department of Statistics, The Chinese University of Hong Kong

Abstract

We propose a likelihood ratio scan method (LRSM) for estimating multiple change-points in piece-wise stationary processes. Using scan statistics reduces the computationally infeasible global multiple change-point estimation problem to a number of single change-point detection problems in various local windows. The computation can be efficiently performed with order $O(n \log n)$. Consistency for the estimated numbers and locations of the change-points are established. Moreover, a procedure is developed for constructing confidence intervals for each of the change-points. Simulation experiments and real data analysis are conducted to illustrate the efficiency of the LRSM.

The swarm intelligence based (SIB) method and its application in statistics

Frederick Kin Hing Phoa

Institute of Statistical Science, Academia Sinica, Taiwan

Abstract

Natural heuristic methods, like the particle swarm optimization and many others, enjoy fast convergence towards optimal solution via a series of inter-particle communication. Such methods are common for the optimization problem in engineering, but few in statistics problem. It is especially difficult to implement in some fields of statistics as the search spaces are mostly discrete, while most natural heuristic methods require continuous search domains. This talk introduces a new method called the Swarm Intelligence Based (SIB) method for optimization in statistics problems, featuring the searches within discrete space. Such fields include experimental designs, community detection, change-point analysis, variable selection, etc. The SIB method is a natural heuristic method that includes the MIX and MOVE operations, which combines target units and selects the best units respectively. This method is advantageous over the traditional particle swarm optimization and many other heuristic approaches in the sense that it is ready for the search of both continuous and discrete domains, and its global best particle is guaranteed to monotonically move towards the optimum. The SIB method is demonstrated in several examples.

Correlation motif model for integrative analyses of genomic data

Ying Ying Wei

Department of Statistics, The Chinese University of Hong Kong

Abstract

In the era of high-throughput technologies, genomic data in public repositories are rapidly growing. To illustrate, to date, more than 1,000,000 samples have been stored in Gene Expression Omnibus and ArrayExpress; meanwhile, over 2,500 ChIP-seq samples are deposited in the ENCODE project and the Sequence Read Archive. This large volume of data provides unprecedented opportunities to improve detection for weak signals by integrating multiple genomic datasets. Here, we propose a scalable correlation motif approach for integrative analysis of multiple high-dimensional genomic datasets. The approach adopts a flexible Bayesian hierarchical mixture model to capture the latent correlation structures embedded in the data and substantially improves signal detection for low-signal-to noise ratio data. The applications are illustrated by differential gene expression detection when the expression datasets have only a small number of replicate samples as well as allele-specific protein-DNA binding detection from ChIP-seq data. Moreover, the proposed model is applicable to heterogeneous data type integration and allows parallel computing.

Multiple comparisons with two controls for ordered categorical responses

Ping Yang, Siu Hung Cheung and Wai Yin Poon

Department of Statistics, The Chinese University of Hong Kong

Abstract

In clinical studies, it is popular that the responses are ordered categorical. To compare efficacy of several treatments with a control with ordinal responses, the normal latent variable model has recently been proposed. This approach conceptualizes the responses as manifestations of an underlying continuous normal variable. In this paper, we extend this idea to develop the multiple comparison method when there are two controls in the clinical trial. The proposed method is constructed such that the familywise type I error rate is being controlled at a pre-specified level. In addition, for a given level of test power, the procedure to evaluate the required sample size is provided. The proposed testing procedure is also illustrated by an example from a clinical study.

Key Words: Ordinal responses; Multiple Comparisons; Multiple Controls; Familywise error rate; Latent normal model.

On the weak convergence and central limit theorem of blurring and nonblurring processes with application to robust location estimation

Ting-Li Chen^{*a}, Hironori Fujisawa^b, Su-Yun Huang^a, Chii-Ruey Hwang^c

a: Institute of Statistical Sciences, Academia Sinica

b: Institute of Statistical Mathematics

c: Institute of Mathematics, Academia Sinica

Abstract

In this talk, I will first present theoretical properties of the blurring and nonblurring processes including their weak convergence to a Brownian bridge-like process and associated Central Limit Theorem. Then I will apply the derived Central Limit Theorem to the estimation of location parameter. I will present simulation studies comparing location estimation based on using blurring and nonblurring processes. The simulation results suggest that location estimation based on the convergence point of blurring process is more robust and often more efficient than that of nonblurring process.

Globally efficient nonparametric inference of average treatment effects by empirical balancing calibration weighting

Phillip Yam

Department of Statistics, The Chinese University of Hong Kong

Abstract

The estimation of average treatment effects based on observational data is extremely important in practice and has been studied by generations of statisticians under different frameworks. Existing globally efficient estimators require non-parametric estimation of a propensity score function, an outcome regression function or both, but their performance can be poor in practical sample sizes. Without explicitly estimating either function, in this talk, I shall consider a wide class calibration weights constructed to attain an exact three-way balance of the moments of observed covariates among the treated, the controls, and the combined group. The wide class includes exponential tilting, empirical likelihood and generalized regression as important special cases, and extends survey calibration estimators to different statistical problems and with important distinctions. Global semiparametric efficiency for the estimation of average treatment effects is established for this general class of calibration estimators. The results show that efficiency can be achieved by solely balancing the covariate distributions without resorting to direct estimation of propensity score or outcome regression. Besides, I shall introduce a consistent estimator of the efficient asymptotic variance, which does not involve additional functional estimation of either the propensity score or the outcome regression functions. The proposed variance estimator outperforms existing estimators that require a direct approximation of the efficient influence function. This is a joint work with Gary Chan (Uni. Washington) and Zheng Zhang (CUHK).

Parameter clustering: encouraging similarities between estimates via euclidean distance regularization

Tso-Jung Yen

Institute of Statistical Science, Academia Sinica

Abstract

In statistical estimation, one important goal is to obtain a model that has better ability in prediction but fewer parameters for interpretation. Such parsimony requirement leads statisticians to develop various techniques for reducing the effective number of parameters in the model. In this paper we propose a penalized estimation method to fulfill this requirement. The method aims to reduce the effective number of parameters by estimating parameters with identical values. It imposes l_2 -norm penalty functions on differences between pairs of the parameters. Under this setting, the method is able to shrink the differences to zero, yielding identical estimates for the parameters. To numerically carry out the method, we first formulate the problem as a constrained optimization problem, and then solve the constrained optimization problem by developing an iterative algorithm based on the alternating direction method of multipliers. Simulation studies show that the method can simultaneously identify the number of effective parameters and deliver collaborative estimates for these parameters. We discuss several applications and a proposal for carrying out this method via distributed optimization.

Keywords: Euclidean distance; Fused lasso; l_2 -norm regularization; Alternating direction method of multipliers; Block splitting algorithms.

Well-posedness of mean-field type forward-backward stochastic differential equations

Alain Bensoussan, Phillip Yam and Zheng Zhang^{*}

Department of Statistics, The Chinese University of Hong Kong

Abstract

Being motivated by a recent pioneer work Carmona, in this talk, we propose a broad class of natural monotonicity conditions under which the unique existence of the solutions to Mean-field type (MFT) forward-backward stochastic differential equations (FBSDE) can be established. Our conditions provided here are consistent with those normally adopted in the traditional FBSDE (without the interference of a mean-field) frameworks, and give a generic explanation on the unique existence of solutions to common MFT-FBSDEs, such as those in the linear-quadratic setting; besides, the conditions are 'optimal' in a certain sense that can elaborate on how their counter-example in Carmona (2013) just fails to ensure its well-posedness. Finally, a comparison theorem is also included.

April 20th, 2015
16:55-17:20

Bayesian detection of embryonic gene expression onset in *C. elegans*

PhD candidate: Jie Hu
Supervisor: Professor Xiaodan Fan

Department of Statistics, The Chinese University of Hong Kong

Abstract

One fundamental question in biology is how a zygote develops into an embryo with different tissues. To approach this goal, large-scale 4D confocal movies of *C.elegans* embryos have been produced by experimental biologists. However, the lack of principled statistical methods for the highly noisy data has hindered the comprehensive analysis of these data sets. We introduced a probabilistic change point model on the cell lineage tree to estimate the embryonic gene expression onset time. A Bayesian approach is used to fit the 4D confocal movies data to the model. Subsequent classification methods are used to decide a model selection threshold and further refine the expression onset time from the branch level to the specific cell time level. Extensive simulations have shown the high accuracy of our method. Its application on real data yielded both previously known results and new findings.

Genome-wide pharmacogenomic study on methadone maintenance treatment identifies SNP rs17180299 and multiple haplotypes on *CYP2B6*, *SPON1*, and *GSG1L* associated with plasma concentrations of methadone R- and S-enantiomers in heroin-dependent patients

Shih-Kai Chu

Institute of Statistical Science, Academia Sinica

*Bioinformatics Program, Taiwan International Graduate Program, Academia Sinica,
Institute of Biomedical Informatics, National Yang-Ming University*

Abstract

Methadone maintenance treatment (MMT) is commonly used for controlling opioid dependence, preventing withdrawal symptoms, and improving the quality of life of heroin-dependent patients. A steady-state plasma concentration of methadone enantiomers, a measure of methadone metabolism, is an index of treatment response and efficacy of MMT. Although the methadone metabolism pathway has been partially revealed, no genome-wide pharmacogenomic study has been performed to identify genetic determinants and characterize genetic mechanisms for the plasma concentrations of methadone R- and S-enantiomers. This study was the first genome-wide pharmacogenomic study to identify genes associated with the plasma concentrations of methadone R- and S-enantiomers and their respective metabolites in a methadone maintenance cohort. After data quality control was ensured, a dataset of 344 heroin-dependent patients in the Han Chinese population of Taiwan who underwent MMT was analyzed. Genome-wide single-locus and haplotype-based association tests were performed to analyze four quantitative traits: the plasma concentrations of methadone R- and S-enantiomers and their respective metabolites. A significant single nucleotide polymorphism (SNP), rs17180299 (raw $p = 2.24 \times 10^{-8}$), was identified, accounting for 9.541% of the variation in the plasma concentration of the methadone R-enantiomers. In addition, 17 haplotypes were identified on *SPON1*, *GSG1L*, and *CYP450* genes associated with the plasma concentration of methadone S-enantiomers. These haplotypes accounted for approximately one-fourth of the variation of the overall S-methadone plasma concentration, in which two significant haplotypes on *CYP2B6* accounted for 10.72% of the variation. A gene expression experiment revealed that *CYP2B6*, *SPON1*, and *GSG1L* can be activated concomitantly through a constitutive androstane receptor (CAR) activation pathway. In conclusion, this study revealed new genes associated with the plasma concentration of methadone, providing insight into the genetic

April 20th, 2015
17:20-17:45

foundation of methadone metabolism. The results can be applied to predict treatment responses and methadone-related deaths for individualized MMTs.

This is joint work with Hsin-Chou Yang, Hsiang-Wei Kuo, Sheng-Chang Wang, Sheng-Wen Liu, Shu Chih Liu, Ing-Kang Ho and Yu-Li Liu.