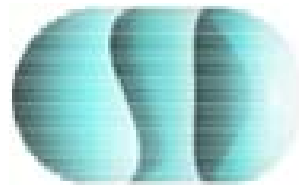


PROCEEDINGS

**The 18th HKBU-CSD Postgraduate Research
Symposium**

PG Day 2015



**Department of Computer Science
Hong Kong Baptist University
May 26th, 2015**

The 18th HKBU-CSD Postgraduate Day Program

May 26 th , 2015, Tuesday (Venue. LMC514)	
Time	Sessions
09:15-09:30	On-site Registration
Opening	
09:30-09:40	Welcome: Prof. P. C Yuen (<i>Head of Department of Computer Science, HKBU</i>)
09:40-09:50	Lu Yang (YMC /LIU Yang, PC)
09:50-10:00	Liu Siqi (PC, LU, YMC /LIU Yang, TAM)
10:00-10:20	Mai Guangcan (PC, YMC /LIU Yang, LU, TAM)
10:20-10:30	Yang Baoyao (PC, YMC /LIU Yang, TAM, LU)
10:30-10:40	Tea Break
10:40-10:50	Shi Qiquan (LU, YMC /LIU Yang, TAM, Yi Wang)
10:50-11:00	Zhou Yichao (YMC /LIU Yang, LU, CHU, Byron)
11:00-11:20	Lou Jian (YMC /LIU Yang, Byron, CHU, LC)
11:20-11:30	Wu Wen (LC, LU, CHU, Byron)
11:30-13:30	Noon Break (All Participants Lunch at RRS Staff Canteen)
13:30-13:40	Zhang Yiqun (YMC /LIU Yang, LU, JX, TAM)
13:40-14:00	Yi Peipei (Byron, JX, YMC /LIU Yang, CHU)
14:00-14:20	Yu Lu (YW, Hai Liu, CHU, JNG, JX)
14:20-14:40	Liu Chengjian (CHU, YW, Hai Liu, JNG)
14:40-15:00	He Jingzhu (YW, CHU, Albert LAM , JNG)
15:00-15:10	Tea Break
15:10-15:20	Xu Cheng (JX, Byron, CHU, Haibo)
15:20-15:40	Chen Qian (JX, Byron, LC, Haibo)
15:40-16:00	Chen Lei (JX, Byron, LC, Haibo)
16:00-16:20	Fan Zhe (Byron, JX, CHU, LC)
16:20-16:40	Gu Fangqing (YMC /LIU Yang, LC, CHU, Byron)
16:40-16:50	Tea Break
16:50-17:00	Wen Xueping (YMC/LIU Yang, LU, William, TAM)
17:00-17:20	Liang Fengfeng (LC, William, JM, JX)
17:20-17:40	Li Chen (William, Byron, JM, LC)
17:40-17:50	Shan Songwei (LC, JM, William, YMC /LIU Yang)
17:50-18:00	Tan Qi (JM, YMC /LIU Yang, William, LC)
18:00-18:20	Bao Qing (William, JM, YMC /LIU Yang)
Best Paper and Best Presentation Awards (To be Presented in the next LunchBite)	
Closing	

Table of abstracts

Section I.....	1
Learning from Imbalanced Data by Hybrid Sampling with Boosting	1
Lu Yang.....	1
Progress on rPPG Based Remote Heart Rate Detection.....	1
Liu Siqi.....	1
Fusing Binary Templates for Multi-biometric Cryptosystems.....	1
Mai Guangcan.....	1
Domain Adaptation for pose estimation in low quality images	2
Yang Baoyao	2
Section II.....	3
Semi-Orthogonal Multilinear PCA with Relaxed Start.....	3
Shi Qiquan.....	3
A Survey of Lip-password Based Speaker Verification	3
Zhou Yichao	3
Incremental Gradient Method with Proximal Average for Empirical Risk Minimization with Composite Regularizer.....	3
Lou Jian.....	3
Implicit Acquisition of User Personality for Augmenting Movie Recommendations.....	4
Wu Wen.....	4
Section III.....	5
Fast On-line Hierarchical Clustering	5
Zhang Yiqun.....	5
Autocomplete Subgraph Query Framework for Graph Databases.....	5
Yi Peipei.....	5
Efficient Channel-Hopping Rendezvous Algorithm Based on Available Channel Set	6
Yu Lu.....	6
DeEC. Reliable Declustered Erasure Codes Placement Algorithm for Large-Scale Storage System	6
Liu Chengjian.....	6
Transmission of Media Files Among Cloud Servers	6
He Jingzhu.....	6
Section IV.....	8
Privacy-Preserving Authentication on Top-k Aggregation Queries.....	8
Xu Cheng	8
Authenticated Online Data Integration Services.....	8
Chen Qian	8
Towards Social-aware Ridesharing Group Query Services	8
Chen Lei.....	8
Keys for Graphs	9
Fan Zhe	9
An Evolutionary Algorithm Based on Decomposition for Multimodal Optimization Problems	9
Gu Fangqing.....	9
Section V.....	10
Feature Selection with Mixed Numerical and Categorical Data.....	10
Wen Xueping.....	10
Mining Spatiotemporal Patterns for Active Disease Surveillance Planning.....	10
Liang Fengfeng	10
Recovering Human Mobility Flow Models and Daily Routine Patterns in a Smart	

Environment.....	10
Li Chen.....	10
Understanding the Impacts of Sociodemographic Profiles on Lung Cancer Risk in Toronto, Canada	11
Shan Songwei.....	11
Inferring Network from Cascade Observation: A Survey.....	12
Tan Qi.....	12
A Component-based Diffusion Model with Structural Diversity for Social Networks	12
Bao Qing	12

Section I.

Learning from Imbalanced Data by Hybrid Sampling with Boosting

Lu Yang

Abstracts. With the increasing availability of large amount of data in a wide range of applications, no matter for industry or academia, it becomes crucial to understand the nature of complex raw data, in order to gain more values from data engineering. Although many problems have been successfully solved by some mature machine learning techniques, the problem of learning from imbalanced data is still one of the challenges in the field of data engineering and machine learning, which attracted growing attention in recent years due to its complexity. In this report, I firstly describe the research gap of the current studies about the selection of sampling methods and the sampling rate problem. Then, a novel method is proposed, which uses the hybrid of undersampling and oversampling to process the data in each boosting iteration. The experiment is conducted on 13 data sets from the UCI data repository and compares the proposed method with other three methods in three evaluation metrics. The experiment shows that the proposed method outperforms the sampling methods which is used individually and indicates that the hybrid sampling is a better choice for the imbalanced data problem.

Progress on rPPG Based Remote Heart Rate Detection

Liu Siqu

Abstracts. Heart rate (HR), or pulse rate, is an important parameter of people's health condition. Recently, some work show that the heart rate can be measured in a remote way. This paper reports my research progress on the topic of remote heart rate detection, including the introduction which contents the summary of related work and preliminary experiment on identifying the problem of existing work. Finally, the conclusion part contents a summary of limitation of recent work and future research plan.

Fusing Binary Templates for Multi-biometric Cryptosystems

Mai Guangan

Abstracts. Biometric cryptosystem has been proven to be one of the promising approaches for template protection. Since most methods in this approach require binary input, to extend it for multiple modalities, binary template fusion is required. This paper addresses the issues of multi-biometrics' performance and security, and proposes a new binary template fusion method which could maximize the fused template discriminability and its entropy by reducing the bits dependency. Three publicly available datasets are used for experiments. Experimental results show that the proposed method outperforms the bit-selection method.

Domain Adaptation for pose estimation in low quality images

Yang Baoyao

Abstracts. This paper addresses a domain adaptation model for estimating human pose in low quality images without label information. Given labeled high quality images (source domain) and unlabelled low quality images (target domain), we propose a Latent Self-Adaptive Support Vector Machine (LSASVM) method to adapt the existing pose estimation model for high quality images to the one for low quality images. To solve the problem of no labeled data in low quality images, we also propose Quality-Guided Transfer (QGT) method to generate the data in low quality images based on the quality information in two domains. And a latent model is utilized to measure the adapted degree of each body part. Although the results of our preliminary experiment are not good enough, the performance of some images increases by our approach.

Section II.

Semi-Orthogonal Multilinear PCA with Relaxed Start

Shi Qiquan

Abstracts. Principal component analysis (PCA) is an unsupervised method for learning low-dimensional features with orthogonal projections. Multilinear PCA methods extend PCA to deal with multidimensional data (tensors) directly via tensor-to-tensor projection or tensor-to-vector projection (TVP). However, under the TVP setting, it is difficult to develop an effective multilinear PCA method with the orthogonality constraint. This paper tackles this problem by proposing a novel Semi-Orthogonal Multilinear PCA (SOMPCA) approach. SO-MPCA learns low-dimensional features directly from tensors via TVP by imposing the orthogonality constraint in only one mode. This novel formulation results in more captured variance and more learned features than full orthogonality. For better generalization, we further introduce a new relaxed start (RS) strategy to get SO-MPCA-RS by fixing the starting projection vectors. Experiments on both face (2D) and gait (3D) data demonstrate that SO-MPCA-RS outperforms other competing algorithms on the whole, and the relaxed start strategy is also effective for other TVP-based PCA methods.

A Survey of Lip-password Based Speaker Verification

Zhou Yichao

Abstracts. The lip-password based visual speaker verification system can be applied in many fields, such as financial security and Human-Computer Interfaces (HCI). This paper does a detailed survey on the current research status about visual speaker verification systems. It focuses on lip feature extraction methods which our lip-password based speaker verification need. Different kinds of lip features are introduced with respective characteristics, as well as the comparison of their performance.

Incremental Gradient Method with Proximal Average for Empirical Risk Minimization with Composite Regularizer

Lou Jian

Abstracts. Empirical risk minimization is widely used in machine learning and data mining to learn the model by optimizing empirical loss averaged from training set. Facing large training set, fast and scalable optimization method is required. Incremental gradient methods, by exploiting the finite sum structure of the loss function, have been shown that converge as fast as full gradient methods while have loss per-iteration cost as stochastic gradient method. Moreover, composite penalties such as group lasso and graph lasso are applied to regularize the ERM model to induce structured sparsity. Such complex nonsmooth regularizer can make the problem hard to optimize. Proximal average is proposed as a better approximation than smoothing method and a more compact alternative approach than ADMM to handle such regularizer. Inspired by both of the recent advances, we propose a new algorithm to efficiently solve ERM problem with composite regularizer.

Implicit Acquisition of User Personality for Augmenting Movie Recommendations

Wu Wen

Abstracts. In recent years, user personality has been recognized as valuable info to build more personalized recommender systems. However, the effort of explicitly acquiring users' personality traits via psychological questionnaire is unavoidably high, which may impede the application of personality based recommenders in real life. In this paper, we focus on deriving users' personality from their implicit behavior in movie domain and hence enabling the generation of recommendations without involving users' efforts. Concretely, we identify a set of behavioral features through experimental validation, and develop inference model based on Gaussian Process to unify these features for determining users' big-five personality traits. We then test the model in a collaborative filtering based recommending framework on two real-life movie datasets, which demonstrates that our implicit personality based recommending algorithm significantly outperforms related methods in terms of both rating prediction and ranking accuracy. The experimental results point out an effective solution to boost the applicability of personality-based recommender systems in online environment.

Section III.

Fast On-line Hierarchical Clustering

Zhang Yiqun

Abstracts. In this report, our research works on on-line hierarchical clustering are introduced. The research target of us is to explore a new approach which can be utilized to construct a hierarchy of streaming data that satisfies the basic characteristics (homogeneity and monotonicity) of standard hierarchical clustering. Not only for the on-line application, it should also be a choice for fast hierarchical clustering. Thus, several on-line schemes are designed and experiments are implemented for them to investigate the effectiveness. Based on the experiment results, the parts remains to be improved in our frameworks are discussed and future plans for improving the performance of it are also stated.

Autocomplete Subgraph Query Framework for Graph Databases

Yi Peipei

Abstracts. Composing queries have been evidently a tedious task. This is particular true to graph queries, as they are typically verbose and prone to typos. It is compounded with the fact that graph schemas can be missing or too loose for helping query formulation. Despite the great success of query formulation aids, in particular, automatic query completion, to the best of our knowledge, auto subgraph query completion has not been investigated. In this prospectus, we propose a novel framework for auto subgraph query completion (called AUTOQ). Given a user query q as input, AUTOQ returns a ranked list of query increments Δq as output, as opposed to answers of q . Further, users may iteratively apply the increments to compose their queries. The main techniques in AUTOQ can be described as follows. First, we design the logical unit for query increments. We propose novel c -prime features of a graph database, which are frequent subgraphs that (i) form larger frequent subgraphs and (ii) is composed by smaller c -prime features in no more than c ways. Second, we represent a query by c -prime features and propose query composition to enlarge the query with a query increment. Third, we propose a novel index called feature DAG (denoted as FDAG) to efficiently locate candidate query increments and rank them according to a given users' intent. We have conducted an extensive experimental evaluation and a shorter usability test. The results show that our proposed techniques reduce query formulation times and the optimizations are effective.

This study forms the foundation of a stream of research of subgraph query feedback for graph databases. For example, AUTOQ can be readily extended to subgraph similarity search. AUTOQ can be extended to make correction on subgraph queries. Finally, the prospectus focuses on data graphs of modest sizes, e.g. molecules and compounds. AUTOQ may be extended to support large graphs such as social networks (Facebook) and co-authorship networks (DBLP).

Efficient Channel-Hopping Rendezvous Algorithm Based on Available Channel Set

Yu Lu

Abstracts. In cognitive radio networks, rendezvous is a fundamental operation by which cognitive users establish a communication link on a commonly-available channel for communications. Some existing rendezvous algorithms provide guaranteed rendezvous (i.e., rendezvous can be achieved within finite time) and they generate channel-hopping (CH) sequences based on the whole channel set. However, some channels may be unavailable (e.g., being used by the licensed users) and these existing algorithms would randomly replace the unavailable channels in the CH sequence. This random replacement is not effective, especially when the number of unavailable channels is large. In this paper, we design a new rendezvous algorithm, called Interleaved Sequences based on *Available Channel set* (ISAC) that attempts rendezvous on the available channels only for faster rendezvous. ISAC constructs an odd subsequence and an even subsequence and interleaves these two subsequences to compose a CH sequence. We prove that ISAC provides guaranteed rendezvous. We derive the upper bounds on the maximum time-to-rendezvous to be $O(m)$ ($m \leq Q$) under the symmetric model and $O(mn)$ ($n \leq Q$) under the asymmetric model, where m and n are the numbers of available channels of two users and Q is the total number of channels. Extensive experiments are conducted to evaluate ISAC.

DeEC. Reliable Declustered Erasure Codes Placement Algorithm for Large-Scale Storage System

Liu Chengjian

Abstracts. Emerged large-scale distributed storage systems employ tens or hundreds of thousands of storage devices in data center or across data center to store petabytes of data. Such systems must adopt reliability mechanisms such as data replication and erasure coding to maximum system reliability. Storage cost for data replication will increase greatly with system growth. This makes erasure coding become a promising reliability mechanism for large-scale distributed storage systems. We proposed DeEC, a declustered data placement algorithm of erasure coding designed for large-scale distributed systems. Because large systems contain many failure domains, our algorithm distributes data and codes across these failure domains to maximum system reliability on multiple layers. The algorithm can accommodate with both traditional block-based storage systems and emerging object-based storage systems by setting proper stripe size for erasure coding.

Transmission of Media Files Among Cloud Servers

He Jingzhu

Abstracts. Cloud servers are promising to support content delivery services. These cloud servers are usually located at different data centers in order to efficiently serve the geographically distributed clients. When new media files are produced, it is necessary to transmit these media files to the distributed cloud servers. Media files have several distinct features. their sizes are usually large, they may have very different sizes, and they may have different popularity. It is desirable that the transmission of the media files can be completed as soon as possible. In this paper, we formulate this transmission problem as an integer linear programming problem where the objective is to minimize the mean completion time. We propose two

methods to solve this problem. In the first method, we decompose the integer linear programming problem into a set of sub-problems and solve them in a parallel manner using the sub-gradient method for the optimal solution. In the second method, we analytically derive some basic features of the problem and exploit these features to design a heuristic algorithm that can find sub-optimal solutions in a short execution time.

Section IV.

Privacy-Preserving Authentication on Top-k Aggregation Queries

Xu Cheng

Abstracts. We study the problem of authenticated top-k aggregation queries over outsources databases. In addition to prove the answer is integrity, we propose the scheme to preserve privacy on the clients side, where the individual records are concealed. We achieve this through an elaborate fusion of *Merkle Grid Tree* (MG-Tree) and *Over-threshold Set Union Authentication* (OTSU). Based on q -Strong Bilinear Diffie-Hellman assumption, we prove the security of our solution. Finally, we experimentally confirm feasibility of proposed scheme.

Authenticated Online Data Integration Services

Chen Qian

Abstracts. Data integration involves combining data from multiple sources and providing users with a unified query interface. Data integrity has been a key problem in online data integration. Although a variety of techniques have been proposed to address the data consistency and reliability issues, there is little work on assuring the integrity of integrated data and the correctness of query results. In this paper, we take the first step to propose authenticated data integration services to ensure data and query integrity even in the presence of an untrusted integration server. We develop a novel authentication code called homomorphic secret sharing seal that can aggregate the inputs from individual sources faithfully by the untrusted server for future query authentication. Based on this, we design two authenticated index structures and authentication schemes for queries on multi-dimensional data. We further study the freshness problem in multisource query authentication and propose several advanced update strategies. Analytical models and empirical results show that our seal design and authentication schemes are efficient and robust under various system settings.

Towards Social-aware Ridesharing Group Query Services

Chen Lei

Abstracts. With the deep penetration of smartphones and geo-locating devices, ridesharing is envisioned as a promising solution to transportation-related problems in metropolitan cities, such as traffic congestion and air pollution. Despite the potential to provide significant societal and environmental benefits, ridesharing has not so far been as popular as expected. Notable barriers include social discomfort and safety concerns when traveling with strangers. To overcome these barriers, in this paper, we propose a new type of Social-aware Ridesharing Group (SaRG) queries which retrieves a group of riders by taking into account their social connections and spatial proximities. While SaRG queries are of practical usefulness, we prove that, however, the SaRG query problem is NP-hard. Thus, we design an efficient algorithm with a set of powerful pruning techniques to tackle this problem. We also present several incremental strategies to accelerate the search speed by reducing repeated computations. Moreover, we propose a novel index

tailored to our problem to further speed up query processing. Experimental results on real datasets show that our proposed algorithms achieve desirable performance.

Keys for Graphs

Fan Zhe

Abstracts. Keys for graphs aim to uniquely identify entities represented by vertices in a graph. We propose a class of keys that are recursively defined in terms of graph patterns, and are interpreted with subgraph isomorphism. Extending conventional keys for relations and XML, these keys find applications in object identification, knowledge fusion and social network reconciliation. As an application, we study the entity matching problem that, given a graph G and a set Σ of keys, is to find all pairs of entities (vertices) in G that are identified by keys in Σ . We show that the problem is intractable, and cannot be parallelized in logarithmic rounds. Nonetheless, we provide two parallel scalable algorithms for entity matching, in MapReduce and a vertex-centric asynchronous model. Using real-life and synthetic data, we experimentally verify the effectiveness and scalability of the algorithms.

An Evolutionary Algorithm Based on Decomposition for Multimodal Optimization Problems

Gu Fangqing

Abstracts. This paper presents a non-parameter method to identify the peaks of the multi-modal optimization problems provided that the peaks are characterized by a smaller objective values than their neighbors and by a relatively large distance from points with smaller objective value. Using the identified peaks as the seeds, we decompose the population into some subpopulations and dynamically allocate the computational effort to different subpopulations. We evaluate the proposed approach on the CEC2015 single objective multi-niche optimization problems. The promising experimental results show its efficacy.

Section V.

Feature Selection with Mixed Numerical and Categorical Data

Wen Xueping

Abstracts. Feature selection is a hot topic in data mining and pattern recognition. Researchers proposed many algorithms to solve it, however there still some key issues remain to be investigate further. A brief introduction and overview of feature selection is presented in this work according to different types of supervised and unsupervised methods. Based on the existing methods, we plan to consider the problem feature selection with mixed data, i.e. the data set consists of both numerical and categorical features, and then a research plan is given.

Mining Spatiotemporal Patterns for Active Disease Surveillance Planning

Liang Fengfeng

Abstracts. Active disease surveillance systems that can timely discover the individual incidences have been proposed recently as a desirable way to prevent the spread of infectious disease. Instead of waiting for the reported cases passively, it searches for the patients in an active way. As the monitoring space is usually very large and available resources are very limited, it is important to design a resource allocation strategy to achieve maximum outputs. However, the complex mechanisms of disease transmission comes with the challenge of fluctuating disease outbreak patterns as well as the irregular incoming flow of imported incidences which makes the task of finding an efficient surveillance strategy a difficult work. In this paper, we present an efficient model for active disease surveillance planning via mining the spatiotemporal patterns of infection risk.

Recovering Human Mobility Flow Models and Daily Routine Patterns in a Smart Environment

Li Chen

Abstracts. With the recent advent of ubiquitous computing and sensor technologies, human mobility data can be acquired for monitoring and analysis purposes, e.g., daily routine identification. Mining mobility data is challenging due to the spatial and temporal variations of the human mobility, even for the same activity. In this paper, we propose a methodology to first summarize indoor human mobility traces as a flow-graph using a probabilistic grammar induction algorithm. Then, we recover salient mobility patterns as subflows in the flow-graph. Thus, such patterns/subflows are expected to be corresponding to the activities that often last for a while, e.g., cooking and cleaning. The weighted kernel k-means algorithm is adopted for the subflow extraction. Finally, we detect the occurrences of the subflows along the mobility traces and obtain their daily routines via the eigen-decomposition. To evaluate the effectiveness of the proposed methodology, we applied it to a publicly available smart home data set containing digital traces of an elder living in a smart house for 219 days. We illustrate how the flow-graphs, subflows and daily routine patterns can be inferred from the mobility data. Our preliminary experimental results show that the

proposed approach can detect subflows which are more specific in terms of their correspondence to activities when compared with a frequent pattern clustering approach.

Understanding the Impacts of Sociodemographic Profiles on Lung Cancer Risk in Toronto, Canada

Shan Songwei

Abstracts. Background Prior research shows significant associations between sociodemographic factors, such as age, gender, socioeconomic status (income, education, etc.), and lung cancer risk characterized by the prevalence of lung cancer in a neighborhood. Most sociodemographic factors are indirect factors, whose impacts on lung cancer risk are manifested through their influences on some direct factors. The correlations between sociodemographic factors and some direct factors of lung cancer, e.g., smoking, physical activity and mental illness, have been widely observed. However, whether there exist significant pathways bridging from sociodemographic factors to lung cancer risk remains to be a question. This study aims to gain some insights into how the sociodemographic factors exert the indirect effects on lung cancer risk, by examining the mediating pathways from sociodemographic factors to lung cancer risk and the their moderating effects on the relationships between direct factors and the development of lung cancer.

Methods In this study we take a close examination of publicly accessible health and census data on sociodemographic profiles, lung cancer risk as well as characteristics on smoking, physical activity and mental illness in Toronto neighborhoods, which were released by Statistics Canada, Toronto Public Health and Toronto Community Health Profiles Partnership. We setup a series of hypotheses related to the mediating pathways from sociodemographic factors to lung cancer risk as well as the moderating effects of them on the relationships between direct factors and the development of lung cancer. Then we test the hypotheses using Partial Least Squares-based Structural Equation Modeling method that is well suited for path analysis and theory exploration.

Results Results show that socioeconomic profile affects lung cancer risk through three pathways. (1) socioeconomic profile has a significant negative effect on smoking ($\beta = -0.131$, $p < 0.1$) and smoking is positively related to lung cancer risk ($\beta = 0.279$, $p < 0.01$); (2) socioeconomic profile has a negative impact on mental illness ($\beta = -0.432$, $p < 0.01$) and mental illness is positively correlated with lung cancer risk ($\beta = 0.502$, $p < 0.01$); (3) socioeconomic profile has a significant positive effect on physical activity ($\beta = 0.444$, $p < 0.01$) and physical activity is negatively related to lung cancer risk ($\beta = -0.137$, $p < 0.05$). In addition, gender profile is inversely associated with smoking ($\beta = 0.340$, $p < 0.01$) and smoking has a positive effect on lung cancer risk ($\beta = 0.279$, $p < 0.01$). Age profile strengthens the effects of smoking and mental illness on lung cancer risk ($\beta = 0.103$, $p < 0.1$; $\beta = 0.111$, $p < 0.1$, respectively).

Conclusions This study has gained a further understanding of the impacts of sociodemographic factors on lung cancer risk by confirming the existence of mediating pathways from sociodemographic factors to the prevalence of lung cancer as well as their moderating effects on the relationships between some direct factors and lung cancer risk. These findings has also provided empirical evidence for health authorities in their efforts on reducing lung cancer risk, for example, by addressing the intermediate factors between socioeconomic profiles and lung cancer risk in the socioeconomic groups that may be affected by the factors of smoking, physical activity and mental illness. As shown in the results, sociodemographic groups may be impacted by different intermediate factors, or by the same factors in varying degrees. Thus, reducing lung cancer risk and mitigating sociodemographic variability in the prevalence of lung cancer should take into account the disparities among various sociodemographic groups and focus on the corresponding impact factors.

Inferring Network from Cascade Observation: A Survey

Tan Qi

Abstracts. Diffusion processes are pervasive in many real-world systems, such as information diffusion in online social networks and viruses spread during the period of epidemic. However, in many cases, such signals propagate over hidden diffusion networks. What we could observe are the time when nodes adopt innovation or become infected. Inferring the network structures from observation data is crucial for understanding the propagation of information and virus. In this paper, we do a survey about the inference algorithms for network reconstruction from cascades data. And we list some challenges in this field for further study.

A Component-based Diffusion Model with Structural Diversity for Social Networks

Bao Qing

Abstracts. Diffusion on social networks refers to the process where opinions are spread via the connected nodes. Given a set of observed information cascades, one can infer the underlying diffusion process for social network analysis. The Independent Cascade Model (IC Model) is a widely adopted diffusion model where a node is assumed to be influenced independently by any one of its neighbors. In reality, there exist different factors governing how a node is influenced by its connected neighbors. For instance, the opinions from the neighbors of the same social group are often similar and thus redundant. In this paper, we extend the IC Model by considering that (1) the information coming from the connected neighbors are similar, and (2) the underlying redundancy can be modeled using a dynamic structural diversity measure of the neighbors. Our proposed model assumes each node to be influenced independently by different communities (or components) of its parent nodes, instead of the individual parent nodes directly. An expectation maximization algorithm is derived to infer the model parameters. We compare the performance of the proposed model with the basic IC Model and its variants using both synthetic data sets and a real-world data set containing news stories and web blogs. Our empirical results show that incorporating the community structure of neighbors and the structural diversity measure into the diffusion model significantly improves the accuracy of the inferred model, at the expense of only a reasonable increase in run-time.