

**IP1****Rapid Learning Systems to Improve Patient Outcomes and Control Health Costs**

In this talk, I will briefly present data mining solutions that analyze millions of patient records, impacting three major areas in healthcare. These include automated quality measurement and decision-support from hospitals EMRs, computer-aided diagnosis systems to identify suspicious lesions on medical images, and rapid learning systems to develop predictive models for personalized medicine. The last is based on a first-of-kind rapid learning system: a Euro-US health IT network spanning cancer centers in 5 nations to learn personalized therapies for lung cancer. The majority of the talk will present case studies that illustrate some of the challenges unique to mining healthcare data, and identify a few promising areas for research. These include the breakdown of traditional assumptions inherent in most mining algorithms, learning from multi-source systems, and the development of predictive models for personalized medicine. We conclude with a glimpse of a more-efficient healthcare future, where treatment decisions are driven by evolving knowledge that is continuously mined from patient records collected in health systems all over the world.

Bharat Rao

SIEMENS Healthcare - Health Services  
bharat.rao@siemens.com

**IP2****Some Assembly Required: Organizing in the 21st Century**

Recent advances in Web Science provide comprehensive digital traces of social actions, interactions, and transactions. These data provide an unprecedented exploratorium to model the socio-technical motivations for creating, maintaining, dissolving, and reconstituting into teams for research, business, or social causes. Using examples from research in team science and massively multiplayer online games, Contractor will argue that Network Science serves as the foundation for the development of social network theories and methods to help advance our ability to understand the emergence of effective teams. More importantly, he will argue that these insights will also enable effective teams by building a new generation of recommender systems that leverage our research insights on the socio-technical motivations for creating ties.

Noshir Contractor

Northwestern University  
nosh@northwestern.edu

**IP3****Cross-Domain Knowledge Transfer in Data Mining**

In data mining, we often encounter situations where we have an insufficient amount of high-quality data in a target domain, but we may have plenty of auxiliary data in related domains. Transfer learning aims to exploit these additional data to improve the learning performance in the target domain. In this talk, I will give an overview on some recent advances in transfer learning for challenging data mining problems. I will present structural transfer-learning solutions under heterogeneous feature representations. I will also survey cross-domain transfer learning solutions in online recommendation, social media and social network mining. I will discuss some current limitations of cross-domain

transfer learning and explore possible future directions.

Qiang Yang

Department of Computer Science,  
Hong Kong University of Science  
qyang@cse.ust.hk

**IP4****Temporal Dynamics and Information Retrieval**

Many digital resources, like the Web, are dynamic and ever-changing collections of information. However, most information retrieval tools developed for interacting with Web content, such as browsers and search engines, focus on a single static snapshot of the information. In this talk, I will present analyses of how Web content changes over time, how people re-visit Web pages over time, and how re-visitation patterns are influenced by changes in user intent and content. These results have implications for many aspects of information retrieval and management including crawling policy, ranking and information extraction algorithms, result presentation, and systems evaluation. I will describe a prototype that supports people in understanding how the information they interact with changes over time, and new retrieval models that incorporate features about the temporal evolution of content to improve core ranking. Finally, I will conclude with an overview of some general challenges that need to be addressed to fully incorporate temporal dynamics in information retrieval and information management systems.

Susan Dumais

Microsoft Research  
sdumais@microsoft.com

**CP1****Sparse Group Lasso: Consistency and Climate Applications**

We address the challenge of designing statistical predictive models for climate data that promote *structured sparsity*. We prove theoretical statistical consistency of estimators with *tree-structured* norm regularizers. We consider one particular model, the *Sparse Group Lasso* (SGL), to construct predictors of land climate using ocean climate variables. Our experimental results demonstrate that the SGL model provides better predictive performance than the current state-of-the-art, remains climatologically interpretable, and is robust in its variable selection.

Soumyadeep Chatterjee, Karsten Steinhaeuser

Department of Computer Science and Engineering  
University of Minnesota, Twin Cities  
chat0129@umn.edu, ksteinha@umn.edu

Arindam Banerjee

University of Minnesota  
banerjee@cs.umn.edu

Snigdhasu Chatterjee

School of Statistics  
University of Minnesota, Twin Cities  
chatterjee@stat.umn.edu

Auroop Ganguly

Northeastern University  
Boston, MA  
a.ganguly@neu.edu

## CP1

**Drought Detection of the Last Century: An Mrf-Based Approach**

Droughts are one of the most damaging climate-related hazards. The late 1960s Sahel drought in Africa and the North American Dust Bowl of the 1930s are two examples of severe droughts that have an impact on society and the environment. Due to the importance of understanding droughts, we consider the problem of their detection based on gridded datasets of precipitation. We formulate the problem as the one of finding the most likely configuration of a Markov Random Field and propose an efficient inference algorithm. We apply this algorithm to the Climate Research Unit precipitation dataset spanning 106 years. The empirical results show that the algorithm successfully identifies the major droughts of the twentieth century in different regions of the world.

Qiang Fu

University of Minnesota, Twin Cities  
qifu@cs.umn.edu

Arindam Banerjee  
University of Minnesota  
banerjee@cs.umn.edu

Stefan Liess, Peter Snyder  
University of Minnesota, Twin Cities  
liess@umn.edu, pksnyder@umn.edu

## CP1

**Toward Data-Driven, Semi-Automatic Inference of Phenomenological Physical Models: Application to Eastern Sahel Rainfall**

First-principles based predictive understanding of complex, dynamic physical phenomena, such as regional precipitation, is quite limited due to the lack of complete phenomenological models underlying their physics. We propose a methodology for *data-driven, semi-automatic* inference of plausible phenomenological models and apply it to derive the model for eastern Sahel rainfall variability. To the best of our knowledge, this is the first model of this phenomenon; several of its components are consistent with the known evidence.

Saurabh V. Pendse

North Carolina State University  
Oak Ridge National Laboratory  
svpendse@ncsu.edu

Isaac Tetteh, Fredrick Semazzi  
North Carolina State University  
itetteh@ncsu.edu, fred\_semazzi@ncsu.edu

Vipin Kumar  
University of Minnesota  
kumar@cs.umn.edu

Nagiza Samatova  
North Carolina State University  
Oak Ridge National Laboratory  
samatova@csc.ncsu.edu

## CP1

**Detecting and Tracking Coordinated Groups in****Dense, Systematically Moving, Crowds**

We address the problem of detecting and tracking clusters of moving objects in very noisy environments. Monitoring a crowded football stadium for small groups of individuals acting suspiciously is an example instance of this problem. In this example the vast majority of individuals are not part of a suspicious group and are considered as noise. Existing spatio-temporal cluster algorithms are only capable of detecting small clusters in extreme noise when the noise objects are moving randomly. In reality, including the example cited, the noise objects move more systematically instead of moving randomly. The members of the suspicious groups attempt to mimic the behaviors of the crowd in order to blend in and avoid detection. This significantly exacerbates the problem of detecting the true clusters. We propose the use of Support Vector Machines (SVMs) to differentiate the true clusters and their members from the systematically moving noise objects. Our technique utilizes the relational history of the moving objects, implicitly tracked in a relationship graph, and a SVM to increase the accuracy of the clustering algorithm. A modified DB-SCAN algorithm is then used to discover clusters of highly related objects from the relationship graph. We evaluate our technique experimentally on several data sets of mobile objects. The experiments show that our technique is able to accurately and efficiently identify groups of suspicious individuals in dense crowds.

James C. Rosswog, Kanad Ghose

Binghamton University  
jim.rosswog@binghamton.edu, ghose@cs.binghamton.edu

## CP1

**Large-Scale Nonparametric Estimation of Vehicle Travel Time Distributions**

Fitting distributions of travel-time in vehicle traffic is an important application of spatio-temporal data mining. While regression methods to forecast the expected travel-time are standard approaches of travel-time prediction, we need to estimate distributions of the travel-time when using state-of-the-art risk-sensitive route recommendation systems. The authors introduce a novel nonparametric density estimator of travel-time for each road or link. The new estimator consists of basis functions modeled as mixtures of gamma or log-normal density functions, a sparse link similarity matrix given as an approximate diffusion kernel on a link connectivity graph, and importance weights for each link. Unlike the existing nonparametric methods that are computationally intensive, the new estimator is stably applicable to large datasets, because the basis functions and the importance weights are globally optimized with a fast convex clustering algorithm. Experimental results using real probe-car datasets show advantages of the new nonparametric estimator over parametric regression methods.

Rikiya Takahashi, Takayuki Osogami, Tetsuro Morimura

IBM Research - Tokyo  
rikiya@jp.ibm.com, osogami@jp.ibm.com,  
tetsuro@jp.ibm.com

## CP2

**The Multi-Set Stream Clustering Problem**

The problem of clustering has been widely studied by the data mining community because of its applications to a wide variety of problems in the context of customer seg-

mentation, electronic commerce and learning. In general, the problem of clustering is generally presented as one of clustering *individual instances* of data records. In many applications, we have a collection of multiple *sets* of records. Each such set is essentially a database of records, and each database may possibly contain a different number of records. It is desirable to cluster these sets on the basis of the *similarity of underlying data distribution*. Thus, this problem may also be understood as that of clustering sets of data sets, as opposed to clustering sets of instances. The problem is especially challenging when the data sets are not available at one time, but are presented in the form of out-of-order and mixed streams, in which the records from different data sets do not arrive in any particular order, but are mixed with one another. In this paper, we present a first approach to the problem with the use of anchor-based summarization. We present experimental results for the effectiveness and efficiency of the approach on a number of real data sets.

Charu C. Aggarwal  
IBM T. J. Watson Research Center  
charu@us.ibm.com

## CP2

### Cluster-Aware Compression with Provable K-Means Preservation

This work rigorously explores the design of cluster-preserving compression schemes for high-dimensional data. We focus on the K-means algorithm and identify conditions under which running the algorithm on the compressed data yields the same clustering outcome as on the original. The compression is performed using single and multi-bit minimum mean square error quantization schemes as well as a given clustering assignment of the original data. We provide theoretical guarantees on post-quantization cluster preservation under certain conditions on the cluster structure, and propose an additional data transformation that can ensure cluster preservation unconditionally; this transformation is invertible and thus induces virtually no distortion on the compressed data. In addition, we provide an efficient scheme for multi-bit allocation, per cluster and data dimension, which enables a trade-off between high compression efficiency and low data distortion. Our experimental studies highlight that the suggested scheme accurately preserved the clusters formed in all cases, while incurring minimal distortion on the data shapes. Our results can find many applications, e.g., in a) clustering, analysis and distribution of massive datasets, where the proposed data compression can boost performance while providing provable guarantees on the clustering result, as well as, in b) cloud computing services, as the optional transformation provides a data-hiding functionality in addition to preserving the K-means clustering outcome.

Nikolaos Freris, Michail Vlachos, Deepak Turaga  
IBM Research  
nif@zurich.ibm.com, michalis0@gmail.com, turaga@us.ibm.com

## CP2

### Supervised Clustering of Label Ranking Data

This paper studies supervised clustering of label ranking data. Potential applications include target marketing, where the goal is to cluster customers in feature space by taking into consideration the assigned, potentially incomplete product preferences. We establish several heuristic

baselines and propose a principled algorithm based on the Plackett-Luce ranking model specifically tailored for this type of clustering. Experimental evaluation on synthetic and real-life data showed that the PL-based method was superior to the baseline approaches.

Mihajlo Grbovic, Nemanja Djuric, Slobodan Vucetic  
Temple University  
mihajlo.grbovic@temple.edu, nemanja.djuric@temple.edu, slobodan.vucetic@temple.edu

## CP2

### Symmetric Nonnegative Matrix Factorization for Graph Clustering

We offer conceptual understanding for the capabilities and shortcomings of nonnegative matrix factorization (NMF) as a clustering method, and propose Symmetric NMF (SymNMF) as a general framework for graph clustering. SymNMF finds a nonnegative symmetric factorization of a matrix containing pairwise similarity values. We then explain why SymNMF captures cluster structures in graph more naturally than spectral clustering. Promising experiment results with Newton-like algorithms are shown using artificial graph data, text data, and image data.

Da Kuang  
Georgia Institute of Technology  
da.kuang@cc.gatech.edu

Chris Ding  
University of Texas at Arlington  
chqding@uta.edu

Haesun Park  
Georgia Institute of Technology  
hpark@cc.gatech.edu

## CP2

### Stratification Based Hierarchical Clustering Over a Deep Web Data Source

This paper focuses on the problem of clustering data from a *hidden* or a deep web data source. A key characteristics of deep web data sources is that data can only be accessed through the limited *query interface* they support. Because the underlying data set cannot be accessed directly, data mining must be performed based on sampling of the datasets. The samples, in turn, can only be obtained by querying the deep web databases with specific inputs. Unlike existing sampling based methods, sampling costs, and not the computation or memory costs, are the dominant consideration in designing the technique for sampling. We have developed a new methodology for addressing the clustering problem on the deep web. Our work includes three new ideas, which are a method for stratifying a deep web data source, an algorithm for hierarchical clustering based on stratified sampling, and a two phase technique for sampling, which includes a representative sampling in the first phase, and sampling focusing on the boundary points between the clusters in the second phase. We have evaluated our approach using two synthetic and one real data set. Our experiments show that each of the three ideas we have introduced leads to significant improvements in accuracy and efficiency of clustering a hidden data source. Specifically, we improve the accuracy of the clusters obtained (measured by average distance to centers) by up to 20% over the existing approach. Compared in another way, our method can achieve the same accuracy with up to 25%

fewer samples, thus reducing the sampling cost.

Tantan Liu, Gagan Agrawal  
The Ohio State University  
liut@cse.ohio-state.edu, agrawal@cse.ohio-state.edu

### CP3

#### Multi-Skill Collaborative Teams Based on Densest Subgraphs

We consider the problem of identifying a team of skilled individuals for collaboration in the presence of a social network with the goal of maximizing collaborative compatibility of the team. The collaborative compatibility is measured as the density of the induced subgraph on selected nodes. We present a 3-approximation algorithm for the single-skill team formation problem and a special case of multiple skills. Our experiments show that these algorithms outperform previous work on several metrics.

Amita Gajewar  
Yahoo! Labs, Yahoo! Inc., Santa Clara, CA  
amitag@yahoo-inc.com

Atish Das Sarma  
Google Research, Google Inc., Mountain View, CA, USA  
dassarma@google.com

### CP3

#### Evaluating Event Credibility on Twitter

Given a set of popular Twitter events (with related users and tweets), we study the problem of automatically assessing credibility of such events. We propose a PageRank-like credibility analysis approach using a multi-typed network of events, tweets, and users. Further, event credibility scores are updated using event graph-based optimization, within each iteration. Experiments on events extracted from millions of tweets show that our methods perform significantly better ( $\sim 86\%$ ) than classifier approach ( $\sim 72\%$ ).

Manish Gupta, Peixiang Zhao  
Univ of Illinois at Urbana-Champaign  
manishg.iitb@gmail.com, pzhao4@illinois.edu

Jiawei Han  
University of Illinois at Urbana-Champaign  
hanj@cs.uiuc.edu

### CP3

#### Har: Hub, Authority and Relevance Scores in Multi-Relational Data for Query Search

In this paper, we propose a framework HAR to study the hub and authority scores of objects, and the relevance scores of relations in multi-relational data for query search. The basic idea of our framework is to consider a random walk in multi-relational data, and study in such random walk, limiting probabilities of relations for relevance scores, and of objects for hub scores and authority scores. The main contribution of this paper is to (i) propose a framework (HAR) that can compute the hub, authority and relevance scores by solving limiting probabilities arising from multi-relational data, and can incorporate input query vectors to handle query-specific search; (ii) show existence and uniqueness of such limiting probabilities so that they can be used for query search effectively; and (iii) develop an it-

erative algorithm to solve a set of tensor (multivariate polynomial) equations to obtain such probabilities. Extensive experimental results on TREC and DBLP data sets suggest that the proposed method is very effective in obtaining relevant results to the querying inputs. In the comparison, we find that the performance of HAR is better than those of HITS, SALSA and TOPHITS.

Xutao Li  
Harbin Institute of Technology, China  
xutaolee08@gmail.com

Michael K. Ng  
Department of Mathematics, Hong Kong Baptist University  
mng@math.hkbu.edu.hk

YUNMING Ye  
Harbin Institute of Technology, China  
yeyunming@hit.edu.cn

### CP3

#### Feature Selection with Linked Data in Social Media

Feature selection is widely used in preparing high-dimensional data for effective data mining. Increasingly popular social media data presents new challenges to feature selection. Social media data consists of (1) traditional high-dimensional, attribute-value data such as posts, tweets, comments, and images, and (2) linked data that describes the relationships between social media users as well as who post the posts, etc. The nature of social media also determines that its data is massive, noisy, and incomplete, which exacerbates the already challenging problem of feature selection. In this paper, we illustrate the differences between attribute-value data and social media data, investigate if linked data can be exploited in a new feature selection framework by taking advantage of social science theories, extensively evaluate the effects of user-user and user-post relationships manifested in linked data on feature selection, and discuss some research issues for future work.

Jiliang Tang  
Arizona State University  
ARIZONA STATE UNIVERISTY  
Jiliang.Tang@asu.edu

Huan Liu  
Arizona State University  
huanliu@asu.edu

### CP3

#### Microscopic Social Influence

Social influences, the phenomena that one individual's actions can induce similar behaviors among his/her friends via their social ties, have been observed prevalingly in socially networked systems. While most existing work focuses on studying general, macro-level influence (e.g., diffusion); equally important is to understand social influence at *microscopic* scales (i.e., at the granularity of single individuals, actions, and time-stamps), which may benefit a range of applications. We propose  $\mu$ SI, a microscopic social-influence model wherein: individuals' actions are modeled as temporary interactions between social network (formed by individuals) and object network (formed by targets of actions); one individual's actions influence his/her friends in a dynamic, network-wise manner (i.e., dependent on

both social and object networks). We develop for  $\mu$ SI a suite of novel inference tools that enable to answer questions of the form: How may an occurred interaction trigger another? More importantly, when and where may a new interaction be observed? We carefully address the computational challenges for inferencing over such semantically rich models by dynamically identifying sub-domains of interest and varying the precision of solutions over different sub-domains. We demonstrate the breadth and generality of  $\mu$ SI using two seemingly disparate applications. In the context of social tagging service, we show how it can help improve the accuracy and freshness of resource recommendation; in the context of mobile phone call service, we show how it can help improve the efficiency of paging operation.

Ting Wang  
IBM T.J. Watson Research Center  
tingwang@us.ibm.com

Mudhakar Srivatsa, Dakshi Agrawal  
IBM Research  
msrivats@us.ibm.com, agrawal@us.ibm.com

Ling Liu  
Georgia Tech  
lingliu@cc.gatech.edu

#### CP4

##### **Heterogeneous Data Fusion Via Space Alignment Using Nonmetric Multidimensional Scaling**

This paper aims to align heterogeneous data spaces into one common space, which makes it possible to analyze relationships between them. We propose a novel graph embedding framework and one such method based on non-metric multidimensional scaling (NMDS). The NMDS criteria using distance rank orders effectively handles both the deformation of original spaces and the alignment between them. Experimental results show its advantages over existing methods using multi-lingual data and document-speech data.

Jaegul Choo  
College of Computing  
Georgia Institute of Technology  
joyfull@cc.gatech.edu

Shawn Bohn, Grant Nakamura, Amanda White  
Pacific Northwest National Laboratory  
shawn.bohn@pnl.gov, grant.nakamura@pnl.gov,  
amanda.white@pnl.gov

Haesun Park  
Georgia Institute of Technology  
hpark@cc.gatech.edu

#### CP4

##### **A Bayesian Nonparametric Joint Factor Model for Learning Shared and Individual Subspaces from Multiple Data Sources**

Joint analysis of multiple data sources is becoming increasingly popular in transfer learning, multi-task learning and cross-domain data mining. One promising approach to model the data jointly is through learning the shared and individual factor subspaces. However, performance of this approach depends on the subspace dimensionalities

and the level of sharing needs to be specified a priori. To this end, we propose a nonparametric joint factor analysis framework for modeling multiple related data sources. Our model utilizes the hierarchical beta process as a nonparametric prior to automatically infer the number of shared and individual factors. For posterior inference, we provide a Gibbs sampling scheme using auxiliary variables. The effectiveness of the proposed framework is validated through its application on two real world problems – transfer learning in text and image retrieval.

Sunil K. Gupta  
Deakin University, Geelong Waurn Ponds Campus,  
Australia  
sunil.gupta@deakin.edu.au

Dinh Phung, Svetha Venkatesh  
Deakin University, Geelong Waurn Ponds Campus  
Victoria, Australia  
dinh.phung@deakin.edu.au,  
svetha.venkatesh@deakin.edu.au

#### CP4

##### **Adaptive Multi-Task Sparse Learning with An Application to Fmri Study**

In this paper, we propose two adaptive multi-task learning methods, adaptive multi-task lasso and adaptive multi-task elastic-net. Both of them can simultaneously conduct model estimation and variable selection across different tasks. Under weak assumptions, we establish the asymptotic oracle property. As a case study, we apply the adaptive multi-task elastic-net to a cognitive science problem, where one wants to discover a compact semantic basis for predicting fMRI images. We show that the adaptive multi-task sparse learning method achieves superior performance and provides some insights into how the brain represents the meanings of words.

Xi Chen  
Carnegie Mellon University  
School of Computer Science  
xichen@cs.cmu.edu

Jingrui He, Rick Lawrence  
IBM T.J. Watson Research Center  
jingruhe@us.ibm.com, ricklawr@us.ibm.com

Jaime Carbonell  
Language Technologies Institute  
Carnegie Mellon University  
jgc@cs.cmu.edu

#### CP4

##### **Heterogeneous Datasets Representation and Learning using Diffusion Maps and Laplacian Pyramids**

The diffusion maps and geometric harmonics provide a method for describing and extending the geometry of high dimensional datasets. These methods suffers from two limitations: First, the assumption that the attributes of the processed dataset are comparable. Second, application of the geometric harmonics requires setting for the correct scale. We propose a method for learning heterogeneous datasets by using diffusion maps for unifying heterogeneous dataset and by replacing the geometric harmonics

with Laplacian pyramid extensions.

Neta Rabin

Yale University  
neta.rabin@yale.edu

#### CP4

##### Learning from Heterogeneous Sources Via Gradient Boosting Consensus

Multiple data sources containing different types of features may be available for a given task. For instance, users' profiles can be used to build recommendation systems. In addition, a model can also use users' historical behaviors and social networks to infer users' interests on related products. We argue that it is desirable to collectively use any available multiple heterogeneous data sources in order to build effective learning models. We call this framework *heterogeneous learning*. In our proposed setting, data sources can include (i) non-overlapping features, (ii) non-overlapping instances, and (iii) multiple networks (i.e. graphs) that connect instances. In this paper, we propose a general optimization framework for heterogeneous learning, and devise a corresponding learning model from gradient boosting. The idea is to minimize the empirical loss with two constraints: (1) There should be consensus among the predictions of overlapping instances (if any) from different data sources; (2) Connected instances in graph datasets may have similar predictions. The objective function is solved by stochastic gradient boosting trees. Furthermore, a weighting strategy is designed to emphasize informative data sources, and deemphasize the noisy ones. We formally prove that the proposed strategy leads to a tighter error bound. This approach consistently outperforms a standard concatenation of data sources on movie rating prediction, number recognition and terrorist attack detection tasks. We observe that the proposed model can improve out-of-sample error rate by as much as 80%.

Xiaoxiao Shi

Computer Department, University of Illinois at Chicago  
xiao.x.shi@gmail.com

Jean-Francois Paiement, David Grangier

AT&T Labs

jpaiement@research.att.com, grangier@research.att.com

Philip Yu

University of Illinois at Chicago  
psyu@cs.uic.edu

#### CP5

##### Marbles: Mining Association Rules Buried in Long Event Sequences

Episodes are sequential patterns that describe events that often occur in the vicinity of each other. In this paper we propose an algorithm that mines association rules between two episodes. We introduce two novel confidence measures for the rules, and aim to limit the output by eliminating redundant rules. We define the class of closed rules, a class that contains all non-redundant output. To make the algorithm efficient, we use pruning steps along the way.

Boris Cule, Nikolaj Tatti, Bart Goethals

University of Antwerp

boris.cule@ua.ac.be,

nikolaj.tatti@ua.ac.be,

bart.goethals@ua.ac.be

#### CP5

##### Class Relevant Pattern Mining in Output-Polynomial Time

The set of so-called relevant patterns is a subset of all itemsets particularly suited for pattern-based classification tasks. So far, no efficient algorithm has been developed for computing the set of relevant patterns: all existing solutions have a worst-case complexity which is exponential in the size of the input and output. In this paper, we investigate new properties of the relevant patterns and develop, thereupon, the first algorithm whose runtime is polynomial in the size of the input and output. As we show in the experimental section, this result is not only of theoretical interest but also of practical importance, often reducing the search space by orders of magnitude.

Henrik Grosskreutz

Fraunhofer IAIS

henrik.grosskreutz@iais.fraunhofer.de

#### CP5

##### Mining Patterns in Networks Using Homomorphism

In recent years many algorithms have been developed for finding patterns in graphs and networks. A disadvantage of these algorithms is that they use subgraph isomorphism to determine the support of a graph pattern; subgraph isomorphism is a well-known NP complete problem. In this paper, we propose an alternative approach which mines tree patterns in networks by using subgraph *homomorphism*. The advantage of homomorphism is that it can be computed in polynomial time, which allows us to develop an algorithm that mines tree patterns in arbitrary graphs in incremental polynomial time. Homomorphism however entails two problems not found when using isomorphism: (1) two patterns of different size can be equivalent; (2) patterns of unbounded size can be frequent. In this paper we formalize these problems and study solutions that easily fit within our algorithm.

Anton Dries

Universitat Pompeu Fabra

anton.dries@upf.edu

Siegfried Nijssen

Katholieke Universiteit Leuven

siegfried.nijssen@cs.kuleuven.be

#### CP5

##### Scalable Induction of Probabilistic Real-Time Automata Using Maximum Frequent Pattern Based Clustering

The paper presents a scalable method for learning probabilistic real-time automata (PRTAs), a new type of model that captures the dynamics of multi-dimensional event logs. In multi-dimensional event logs, events are described by several features instead of only one symbol. Moreover, it is not clear up front which events occur in an event log. The learning method to find a PRTA that models such an event log is based on the state merging of a prefix tree acceptor, which is guided by a clustering to determine the states of the automaton. To make the overall approach scalable, an online clustering method based on maximum frequent patterns (MFPs) is used. The approach is evaluated on a synthetic, a biological and a medical data set. The results show that the induction of automata using

MFP-based clustering gives easy to understand and stable automata, but most importantly, makes it scalable to large data sets.

Jana Schmidt  
 Institut fuer Informatik, TU Muenchen  
 jana.schmidt@in.tum.de

Sonja Ansorge  
 TU Muenchen  
 sonja.ansorge@gmx.de

Stefan Kramer  
 Johannes Gutenberg-Universitaet Mainz  
 kramer@informatik.uni-mainz.de

### CP5 **Slim: Directly Mining Descriptive Patterns**

Mining small, useful, and high-quality sets of patterns has recently become an important topic in data mining. The standard approach is to first mine many candidates, and then to select a good subset. However, the pattern explosion generates such enormous amounts of candidates that by post-processing it is virtually impossible to analyse dense or large databases in any detail. We introduce SLIM, an any-time algorithm for mining high-quality sets of itemsets directly from data. We use MDL to identify the best set of itemsets as that set that describes the data best. To approximate this optimum, we iteratively use the current solution to determine what itemset would provide most gain—estimating quality using an accurate heuristic. Without requiring a pre-mined candidate collection, SLIM is parameter-free in both theory and practice. Experiments show we mine high-quality pattern sets; while evaluating orders-of-magnitude fewer candidates than our closest competitor, KRIMP, we obtain much better compression ratios—closely approximating the locally-optimal strategy. Classification experiments independently verify we characterise data very well.

Koen Smets, Jilles Vreeken  
 Universiteit Antwerpen  
 koen.smets@ua.ac.be, jilles.vreeken@ua.ac.be

### CP6 **Beam Methods for the Profile Hidden Markov Model**

The Profile Hidden Markov Model (PHMM) is commonly used to represent biological sequences. We present a method for transforming the Profile HMM into an equivalent standard HMM where each transition is associated with a single emission. Using this transformation, we develop a beam method, which includes a novel variational adaptation of the infinite-HMM beam sampling technique, to create a fast inference algorithm. We evaluate our algorithm on both synthetic data and protein sequence datasets, showing that our beam method can lead to considerable improvements in runtime while maintaining the model's ability to concisely represent sequences.

Samuel J. Blasiak, Huzefa Rangwala, Kathryn Laskey  
 George Mason University  
 sblasiak@gmu.edu, rangwala@cs.gmu.edu,  
 klaskey@gmu.edu

### CP6 **Optimal Distance Estimation Between Compressed Data Series**

Most real-world data contain repeated or periodic patterns. This suggests that they can be effectively represented and compressed using only a few coefficients of an appropriate complete orthogonal basis (e.g., Fourier, Wavelets, Karhunen Loeve expansion or Principal Components). In the face of ever increasing data repositories and given that most mining operations are distance-based, it is vital to perform accurate distance estimation directly on the compressed data. However, distance estimation when the data are represented using different sets of coefficients is still a largely unexplored area. This work studies the optimization problems related to obtaining the tightest lower/upper bound on the distance based on the available information. In particular, we consider the problem where a distinct set of coefficients is maintained for each sequence, and the  $L_2$ -norm of the compression error is recorded. We establish the properties of optimal solutions, and leverage the theoretical analysis to develop a fast algorithm to obtain an exact solution to the problem. The suggested solution provides the tightest provable estimation of the  $L_2$ -norm or the correlation, and executes at least two order of magnitudes faster than a numerical solution based on convex optimization. The contributions of this work extend beyond the purview of periodic data, as our methods are applicable to any sequential or high-dimensional data as well as to any orthogonal data transformation used for the underlying data compression scheme.

Nikolaos Freris, Michail Vlachos  
 IBM Research  
 nif@zurich.ibm.com, michalis0@gmail.com

Serdar Kozat  
 Koc University  
 Istanbul, Turkey  
 skozat@ku.edu.tr

### CP6 **Transformation Based Ensembles for Time Series Classification**

Until recently, the vast majority of data mining time series classification (TSC) research has focused on alternative distance measures for 1-Nearest Neighbour (1-NN) classifiers based on either the raw data, or on compressions or smoothing of the raw data. Despite the extensive evidence in favour of 1-NN classifiers with Euclidean or Dynamic Time Warping distance, there has also been a flurry of recent research publications proposing classification algorithms for TSC. Generally, these classifiers describe different ways of incorporating summary measures in the time domain into more complex classifiers. Our hypothesis is that the easiest way to gain improvement on TSC problems is to simply transform into an alternative data space where the discriminatory features are more easily detected. To test our hypothesis, we perform a range of benchmarking experiments in the time domain, before evaluating nearest neighbour classifiers on data transformed into the power spectrum, the autocorrelation function, and the principal component space. We demonstrate that on some problems there is dramatic improvement in the accuracy of classifiers built on the transformed data over classifiers built in the time domain, but that there is also a wide variance in accuracy for a particular classifier built on different data transforms. To overcome this variability, we propose a simple transformation based ensemble, then demonstrate that

it improves performance and reduces the variability of classifiers built in the time domain only. Our advice to a practitioner with a real world TSC problem is to try transforms before developing a complex classifier; it is the easiest way to get a potentially large increase in accuracy, and may provide further insights into the underlying relationships that characterise the problem.

Anthony Bagnall, Luke Davis, Jon Hills, Jason Lines  
University of East Anglia  
Anthony.Bagnall@uea.ac.uk, luke.davis@uea.ac.uk,  
j.hills@uea.ac.uk, j.lines@uea.ac.uk

## CP6

### Mining Compressing Sequential Patterns

Compression based pattern mining has been successfully applied to many data mining tasks. We propose an approach based on the minimum description length principle to extract sequential patterns that compress a database of sequences well. We show that mining compressing patterns is NP-Hard and belongs to the class of inapproximable problems. We propose two heuristic algorithms to mining compressing patterns. The first uses a two-phase approach similar to Krimp for itemset data. To overcome performance with the required candidate generation we propose GoKrimp, an effective greedy algorithm that directly mines compressing patterns. We conduct an empirical study on six real-life datasets to compare the proposed algorithms by run time, compressibility, and classification accuracy using the patterns found as features for SVM classifiers.

Hoang Thanh Lam  
TU Eindhoven  
t.l.hoang@tue.nl

Fabian Moerchen, Dmitriy Fradkin  
Siemens Corporation, Corporate Research  
fabian.moerchen@siemens.com,  
dmitriy.fradkin@siemens.com

Toon Calders  
TU Eindhoven  
t.calders@tue.nl

## CP6

### Simplex Distributions for Embedding Data Matrices over Time

Early stress recognition is of great relevance in precision plant protection. Pre-symptomatic water stress detection is of particular interest, ultimately helping to meet the challenge of “How to feed a hungry world?”. Due to the climate change, this is of considerable political and public interest. Due to its large-scale and temporal nature, e.g., when monitoring plants using hyperspectral imaging, and the demand of physical meaning of the results, it presents unique computational problems in scale and interpretability. However, big data matrices over time also arise in several other real-life applications such as stock market monitoring where a business sector is characterized by the ups and downs of each of its companies per year or topic monitoring of document collections. Therefore, we consider the general problem of embedding data matrices into Euclidean space over time without making any assumption on the generating distribution of each matrix. To do so, we represent all data samples by means of convex combinations of only few extreme ones computable in linear time. On the simplex spanned by the extremes, there are then natu-

ral candidates for distributions inducing distances between and in turn embeddings of the data matrices. We evaluate our method across several domains, including synthetic, text, and financial data as well as a large-scale dataset on water stress detection in plants with more than 3 billion matrix entries. The results demonstrate that the embeddings are meaningful and fast to compute. The stress detection results were validated by a domain expert and conform to existing plant physiological knowledge.

Kristian Kersting  
Fraunhofer IAIS  
kristian.kersting@iais.fraunhofer.de

Mirwaes Wahabzada  
Fraunhofer IAIS  
Sankt Augustin, Germany  
mirwaes.wahabzada@iais.fraunhofer.de

Christoph Roemer  
Institute of Geodesy and Geoinformation  
University of Bonn, Germany  
roemer@igg.uni-bonn.de

Christian Thureau  
Fraunhofer IAIS, Sankt Augustin, Germany  
christian.thureau@iais.fraunhofer.de

Agim Ballvora  
Institute of Crop Science and Resource Conservation  
Plant Breeding, University of Bonn, Germany  
ballvora@uni-bonn.de

Uwe Rascher  
FZ Jülich  
u.rascher@fz-juelich.de

Jens Leon  
Institute of Crop Science and Resource Conservation  
Plant Breeding, University of Bonn, Germany  
j.leon@uni-bonn.de

Christian Bauckhage  
Fraunhofer IAIS  
christian.bauckhage@iais.fraunhofer.de

Lutz Pluemer  
Institute of Geodesy and Geoinformation  
University of Bonn, Germany  
pluemer@igg.uni-bonn.de

## CP7

### Bayesian Supervised Multilabel Learning with Coupled Embedding and Classification

We introduce a novel *Bayesian supervised multilabel learning* method that combines linear dimensionality reduction with linear binary classification. We present a deterministic variational approximation approach to learn the proposed probabilistic model for multilabel classification. Experiments show that the proposed approach achieves good performance values in terms of hamming loss, macro  $F_1$ , and micro  $F_1$  on held-out test data. The low-dimensional embeddings obtained by our method are also very useful for exploratory data analysis.

Mehmet Gönen  
Department of Information and Computer Science



Aalto University School of Science  
mehmet.gonen@aalto.fi

## CP7

### Multi-Objective Multi-Label Classification

Multi-label classification refers to the task of predicting potentially multiple labels for a given instance. Conventional multi-label classification approaches focus on the single objective setting, where the learning algorithm optimizes over a single performance criterion (e.g. *Ranking Loss*) or a heuristic function. The basic assumption is that the optimization over one single objective can improve the overall performance of multi-label classification and meet the requirements of various applications. However, in many real applications, an optimal multi-label classifier may need to consider the tradeoffs among multiple conflicting objectives, such as minimizing *Hamming Loss* and maximizing *Micro F1*. In this paper, we study the problem of *multi-objective multi-label classification* and propose a novel solution (called MOML) to optimize over multiple objectives simultaneously. Note that optimization objectives may be conflicting, thus one cannot identify a single solution that is optimal on all objectives. Our MOML algorithm finds a set of *non-dominated solutions* which are optimal according to the different tradeoffs of the multiple objectives. So users can flexibly construct various combined predictive models from the solution set, which helps to provide more meaningful classification results in different application scenarios. Empirical studies on real-world tasks demonstrate that the MOML can effectively boost the overall performance of multi-label classification, not limiting to the optimization objectives.

Chuan Shi  
Beijing University of Posts and Telecommunications  
shichuan@bupt.edu.cn

Xiangnan Kong, Philip Yu  
University of Illinois at Chicago  
kongxn@gmail.com, psyu@cs.uic.edu

Bai Wang  
Beijing University of Posts and Telecommunications  
wangbai@bupt.edu.cn

## CP7

### A Distributed Kernel Summation Framework for General-Dimension Machine Learning

Kernel summations are a ubiquitous key computational bottleneck in many data analysis methods. We provide the first distributed implementation of kernel summation framework that can utilize: 1) various types of deterministic and probabilistic approximations; 2) any multi-dimensional binary utilizing both distributed and shared memory parallelism; 3) a dynamic load balancing scheme. We show scalability results for kernel density estimation on a subset of the Sloan Digital Sky Survey Data up to 6,144 cores.

Dongryeol Lee, Richard Vuduc, Alexander Gray  
Georgia Institute of Technology  
dongryel@cc.gatech.edu, richie@cc.gatech.edu, agray@cc.gatech.edu

## CP7

### Subtree Replacement in Decision Tree Simplifica-

## tion

The current availability of efficient algorithms for decision tree induction makes intricate post-processing techniques worth to be investigated both for efficiency and effectiveness. We study the simplification operator of subtree replacement, also known as *grafting*, originally implemented in the C4.5 system. We present a parametric bottom-up algorithm integrating grafting with the standard pruning operator, and analyze its complexity in terms of the number of nodes visited. Immediate instances of the parametric algorithm include extensions of error based, reduced error, minimum error, and pessimistic error pruning. Experimental results show that the computational cost of grafting is paid off by statistically significant smaller trees without accuracy loss.

Salvatore Ruggieri  
Dipartimento di Informatica, Università di Pisa  
ruggieri@di.unipi.it

## CP7

### Kernelized Probabilistic Matrix Factorization: Exploiting Graphs and Side Information

We propose a new matrix completion algorithm—Kernelized Probabilistic Matrix Factorization (KPMF), which effectively incorporates external side information into the matrix factorization process. Unlike Probabilistic Matrix Factorization (PMF), which assumes an independent latent vector for each row (and each column) with Gaussian priors, KPMF works with latent vectors spanning all rows (and columns) with Gaussian Process (GP) priors. Hence, KPMF explicitly captures the underlying (nonlinear) covariance structures across rows and columns. This crucial difference greatly boosts the performance of KPMF when appropriate side information, e.g., users' social network in recommender systems, is incorporated. We demonstrate the efficacy of KPMF through two different applications: 1) recommender systems and 2) image restoration.

Tinghui Zhou  
Carnegie Mellon University  
tinghuiz@cmu.edu

Hanhui Shan  
Department of Computer Science and Engineering  
University of Minnesota, Twin Cities  
shan@cs.umn.edu

Arindam Banerjee  
University of Minnesota  
banerjee@cs.umn.edu

Guillermo Sapiro  
University of Minnesota  
Dept Electrical & Computer Engineering  
guille@umn.edu

## CP8

### On Dynamic Link Inference in Heterogeneous Networks

Network and linked data have become quite prevalent in recent years because of the ubiquity of the web and social media applications, which are inherently network oriented. Such networks are massive, dynamic, contain a lot of content, and may evolve over time in terms of the underlying

structure. In this paper, we will study the problem of dynamic link inference in *temporal* and *heterogeneous* information networks. The problem of dynamic link inference is extremely challenging in massive and heterogeneous information network because of the challenges associated with the dynamic nature of the network, and the different types of nodes and attributes in it. Both the topology and type information need to be used effectively for the link inference process. We propose an effective two-level scheme which makes efficient macro- and micro-decisions for combining structure and content in a *dynamic and time-sensitive* way. The time-sensitive nature of the links is leveraged in order to perform effective link prediction. We illustrate the effectiveness of our technique over a number of real data sets.

Charu C. Aggarwal  
IBM T. J. Watson Research Center  
charu@us.ibm.com

Yan Xie, Philip Yu  
University of Illinois at Chicago  
yxie8@uic.edu, psyu@cs.uic.edu

### CP8

#### Parameter-Free Identification of Cohesive Subgroups in Large Attributed Graphs

Given a graph with node attributes, how can we find meaningful patterns such as clusters, bridges, and outliers? Attributed graphs appear in real world in the form of social networks with user interests, gene interaction networks with gene expression information, phone call networks with customer demographics, and many others. In effect, we want to group the nodes into clusters with similar connectivity and homogeneous attributes. Most existing graph clustering algorithms either consider only the connectivity structure of the graph and ignore the node attributes, or require several user-defined parameters such as the number of clusters. We propose PICS, a novel, parameter-free method for mining *attributed graphs*. Two key advantages of our method are that (1) it requires *no* user-specified parameters such as the number of clusters and similarity functions, and (2) its running time scales *linearly* with total graph and attribute size. Our experiments show that PICS reveals meaningful and insightful patterns and outliers in both synthetic and real data sets, including call networks, political books, political blogs, and collections from Twitter and YouTube which have more than 70K nodes and 30K attributes.

Leman Akoglu  
Carnegie Mellon University  
lakoglu@cs.cmu.edu

Hanghang Tong  
IBM T.J. Watson  
htong@us.ibm.com

Brendan Meeder, Christos Faloutsos  
Carnegie Mellon University  
bmeeder@cs.cmu.edu, christos@cs.cmu.edu

### CP8

#### Structural Analysis in Multi-Relational Social Networks

Modern social networks often consist of multiple relations among individuals. Understanding the structure of such

multi-relational network is essential. In sociology, one way of structural analysis is to identify different positions and roles using blockmodels. In this paper, we generalize stochastic blockmodels to *Generalized Stochastic Blockmodels* (GSBM) for performing positional and role analysis on multi-relational networks. Our GSBM generalizes many different kinds of *Multivariate Probability Distribution Function* (MVPDF) to model different kinds of multi-relational networks. In particular, we propose to use *multivariate Poisson distribution* for multi-relational social networks.

Bing Tian Dai, Freddy Chua, Ee-Peng Lim  
Singapore Management University  
btdai@smu.edu.sg, freddy.chua.2009@phdis.smu.edu.sg, eplim@smu.edu.sg

### CP8

#### Influence Blocking Maximization in Social Networks under the Competitive Linear Threshold Model

In many real-world situations, different and often opposite opinions, innovations, or products are competing with one another for their social influence in a networked society. In this paper, we study competitive influence propagation in social networks under the competitive linear threshold (CLT) model, an extension to the classic linear threshold model. Under the CLT model, we focus on the problem that one entity tries to block the influence propagation of its competing entity as much as possible by strategically selecting a number of seed nodes that could initiate its own influence propagation. We call this problem the influence blocking maximization (IBM) problem. We prove that the objective function of IBM in the CLT model is submodular, and thus a greedy algorithm could achieve  $1 - 1/e$  approximation ratio. However, the greedy algorithm requires Monte-Carlo simulations of competitive influence propagation, which makes the algorithm not efficient. We design an efficient algorithm CLDAG, which utilizes the properties of the CLT model, to address this issue. We conduct extensive simulations of CLDAG, the greedy algorithm, and other baseline algorithms on real-world and synthetic datasets. Our results show that CLDAG is able to provide best accuracy in par with the greedy algorithm and often better than other algorithms, while it is two orders of magnitude faster than the greedy algorithm.

Xinran He, Guojie Song  
Peking University  
xinranhe@pku.edu.cn, gjsong@pku.edu.cn

Wei Chen  
Microsoft Research Asia  
weic@microsoft.com

Qingye Jiang  
Columbia University  
qj2116@columbia.edu

### CP8

#### A Framework for the Evaluation and Management of Network Centrality

Network-analysis literature is rich in node-centrality measures that quantify the centrality of a node as a function of the (shortest) paths of the network that go through it. Existing work focuses on defining instances of such measures and designing algorithms for the specific combinato-

rial problems that arise for each instance. In this work, we propose a unifying definition of centrality that subsumes all path-counting based centrality definitions: e.g., stress, betweenness or paths centrality. We also define a generic algorithm for computing this generalized centrality measure for every node and every group of nodes in the network. Next, we define two optimization problems:  $k$ -GROUP CENTRALITY MAXIMIZATION and  $k$ -EDGE CENTRALITY BOOSTING. In the former, the task is to identify the subset of  $k$  nodes that have the largest group centrality. In the latter, the goal is to identify up to  $k$  edges to add to the network so that the centrality of a node is maximized. We show that both of these problems can be solved efficiently for arbitrary centrality definitions using our general framework. In a thorough experimental evaluation we show the practical utility of our framework and the efficacy of our algorithms.

Vatche Ishakian, Dora Erdos, Evimaria Terzi, Azer Bestavros  
Boston University  
visahak@bu.edu, edori@cs.bu.edu, evimaria@cs.bu.edu, best@cs.bu.edu

### CP9

#### Feature Selection “Tomography” — Illustrating That Optimal Feature Filtering Is Hopelessly Un-generalizable

Feature filtering methods are used in high-dimensional domains to quickly score each feature independently. We provide a new empirical method to reveal the *feature preference surface* for a given situation. This visualization reveals new insights for feature filtering: (a) Existing functions do not match the surfaces we revealed. (b) The shape of the surfaces varies and depends on more factors than have been studied at once in the existing literature in feature filtering.

George Forman  
HP Labs  
ghforman@hpl.hp.com

### CP9

#### A Bayesian Markov-Switching Model for Sparse Dynamic Network Estimation

Inferring Dynamic Bayesian Networks (DBNs) from multivariate time series data is a key step towards the understanding of complex systems as it reveals important dependency relationship underlying such systems. Most of the traditional approaches assume a “static” DBN. Yet in many relevant applications, such as those arising in biology and social sciences, the dependency structures may vary over time. In this paper, we introduce a sparse Markov-switching vector autoregressive model to capture the structural changes in the dependency relationships over time. Our approach accounts for such structural changes via a set of latent state variables, which are modeled by a discrete-time discrete-state Markov process. Assuming that the underlying structures are sparse, we estimate the networks at each state through the hierarchical Bayesian group Lasso, so as to efficiently capture dependencies with lags greater than one time unit. For computation, we develop an efficient algorithm based on the Expectation-Maximization method. We demonstrate the strength of our approach through simulation studies and a real data set concerning

climate change.

Huijing Jiang  
IBM T.J. Watson Research Center  
huijiang@us.ibm.com

Aurelie Lozano  
IBM Research  
T. J. Watson Research Center  
aclozano@us.ibm.com

Fei Liu  
IBM T.J. Watson Research Center  
feiliu@us.ibm.com

### CP9

#### Feature Selection over Distributed Data Streams through Optimization

Monitoring data streams in a distributed system has attracted considerable interest in recent years. The task of feature selection (e.g., by monitoring the information gain of various features) requires a very high communication overhead when addressed using straightforward centralized algorithms. While most of the existing algorithms deal with monitoring simple aggregated values such as frequency of occurrence of stream items, motivated by recent contributions based on geometric ideas we present an alternative approach. The proposed approach enables monitoring values of an arbitrary threshold function over distributed data streams through constraints applied separately on each stream. We report numerical experiments on a real-world data that detect instances where communication between nodes is required, and compare the approach and the results to those recently reported in the literature.

Jacob Kogan  
umbc  
kogan@umbc.edu

### CP9

#### Sampling Strategies to Evaluate the Performance of Unknown Predictors

The focus of this paper is on how to select a small sample of examples for labeling that can help us to evaluate many different classification models unknown at the time of sampling. We are particularly interested in studying the sampling strategies for problems in which the prevalence of the two classes is highly biased toward one of the classes. The evaluation measures of interest we want to estimate as accurately as possible are those obtained from the contingency table. We provide a careful theoretical analysis on sensitivity, specificity, and precision and show how sampling strategies should be adapted to the rate of skewness in data in order to effectively compute the three aforementioned evaluation measures.

Hamed Valizadegan, Saeed Amizadeh, Milos Hauskrecht  
University of Pittsburgh  
hamed@cs.pitt.edu, saeed@cs.pitt.edu, milos@cs.pitt.edu

### CP9

#### Learning Hierarchical Relationships among Partially Ordered Objects with Heterogeneous At-

## tributes and Links

Objects linking with many other objects in an information network may imply various semantic relationships. In this work we study a generic form of relationship along which objects can form a tree-like structure, a pervasive structure in various domains. We formalize the problem of uncovering hierarchical relationships in a supervised setting. We propose a discriminative undirected graphical model which integrates a wide range of features and rules by defining potential functions with simple forms.

Chi Wang

University of Illinois at Urbana-Champaign  
chiwang1@illinois.edu

Jiawei Han

UIUC

hanj@illinois.edu

Qi Li, Xiang Li, Wen-Pin Lin, Heng Ji

City University of New York

liqiearth@gmail.com,

jackieiu729@gmail.com,

danniellin@gmail.com, hengjicuny@gmail.com

## CP10

### Transfer Learning of Distance Metrics by Cross-Domain Metric Sampling Across Heterogeneous Spaces

In this paper, we examine a new angle to the transfer learning problem, where we examine the problem of distance function learning. Specifically, we focus on the problem of how our knowledge of distance functions in one domain can be transferred to a new domain. A good semantic understanding of the feature space is critical in providing the domain specific understanding for setting up good distance functions. Unfortunately, not all domains have feature representations which are equally interpretable. For example, in some domains such as text, the semantics of the feature representation are clear, as a result of which it is easy for a domain expert to set up distance functions for specific kinds of semantics. In the case of image data, the features are semantically harder to interpret, and it is harder to set up distance functions, especially for particular semantic criteria. In this paper, we focus on the problem of transfer learning as a way to close the semantic gap between different domains, and show how to use correspondence information between two domains in order to set up distance functions for the semantically more challenging domain.

Guo-Jun Qi

Beckman Institute

University of Illinois at Urbana-Champaign  
qi4@illinois.edu

Charu C. Aggarwal

IBM T. J. Watson Research Center  
charu@us.ibm.com

Thomas Huang

University of Illinois at Urbana-Champaign  
huang@ifp.uiuc.edu

## CP10

### Transfer Topic Modeling with Ease and Scalability

The increasing volume of *short* texts generated on social media sites, such as Twitter or Facebook, creates a

great demand for effective and efficient topic modeling approaches. While latent Dirichlet allocation (LDA) can be applied, it is not optimal due to its weakness in handling short texts with fast-changing topics and scalability concerns. In this paper, we propose a transfer learning approach that utilizes abundant labeled documents from other domains (such as Yahoo! News or Wikipedia) to improve topic modeling, with better model fitting and result interpretation. Specifically, we develop *Transfer Hierarchical* LDA (thLDA) model, which incorporates the label information from other domains via informative priors. In addition, we develop a parallel implementation of our model for large-scale applications. We demonstrate the effectiveness of our thLDA model on both a microblogging dataset and standard text collections including AP and RCV1 datasets.

Jeon-Hyung Kang, Jun Ma, Yan Liu

University of Southern California

jeonhyuk@usc.edu, junma@usc.edu, yanliu@usc.edu

## CP10

### Dual Transfer Learning

In this paper, we propose a novel approach, Dual Transfer Learning (DTL), which simultaneously learns the marginal and conditional distributions, and exploits the duality between them in a principled way. The key idea behind DTL is that learning one distribution can help to learn the other. This duality property leads to mutual reinforcement when adapting both distributions across domains to transfer knowledge. Experiments demonstrate the effectiveness of our proposed approach.

Mingsheng Long

Department of Computer Science and Technology

Tsinghua University

longmingsheng@gmail.com

Jianmin Wang, Guiguang Ding

School of Software

Tsinghua University

jimwang@tsinghua.edu.cn, dinggg@tsinghua.edu.cn

Wei Cheng

Department of Computer Science

University of North Carolina at Chapel Hill

weicheng@cs.unc.edu

Xiang Zhang

Department of Electrical Engineering and Computer Science

Case Western Reserve University

xiang.zhang@case.edu

Wei Wang

Department of Computer Science

University of North Carolina at Chapel Hill

weiwang@cs.unc.edu

## CP10

### Transfer Significant Subgraphs Across Graph Databases

A key step of graph classification is to identify informative subgraphs that encode label information. For instance, in drug efficacy prediction, the drugs (chemical compounds) effective against the same disease usually contain similar

chemical-subgraphs effective to control the disease. Then, one can use such chemical subgraphs to identify effective drugs. We call these subgraphs *significant subgraphs*. In this paper, the aim is to utilize the significant subgraphs from related graph datasets to help label graphs of the target dataset. For example, we utilize the breast cancer drug data, and transfer the anti-cancer subgraphs to help label another set of drug data against lung cancer. To do so, we propose a Bayesian-based transfer learning model. The key idea is to first evaluate the similarity between the target and source datasets by estimating the degree they share on their significant subgraphs. This dataset similarity is then used to judiciously select significant subgraphs from similar (related) datasets to the target dataset. An optimization problem is devised to maximize the likelihood that the selected subgraphs are significant in the target dataset. The objective function is further proven to have the antimotone property which can help prune the search space significantly. Sixteen sets of experiments show that the proposed algorithm can effectively reduce the error rates by as much as 40%. More importantly, it is 10 times faster than the comparison models, which include unsupervised and supervised significant subgraph mining algorithms.

Xiaoxiao Shi

Computer Department, University of Illinois at Chicago  
xiao.x.shi@gmail.com

Xiangnan Kong, Philip Yu  
University of Illinois at Chicago  
kongxn@gmail.com, psyu@cs.uic.edu

#### CP11

##### **Sor: Scalable Orthogonal Regression for Low-Redundancy Feature Selection and Its Healthcare Applications**

As more clinical information with increasing diversity become available for analysis, a large number of features can be constructed and leveraged for predictive modeling. Feature selection is a classic analytic component that faces new challenges due to the new applications: How to handle a diverse set of high dimensional features? How to select features with high predictive power, but low redundant information? How to design methods that can select globally optimal features with theoretical guarantee? How to incorporate and extend existing knowledge driven approach? In this paper, we present Scalable Orthogonal Regression (SOR), an optimization-based feature selection method with the following novelties: 1) Scalability: SOR achieves nearly linear scale-up with respect to the number of input features and the number of samples; 2) Optimality: SOR is formulated as an alternative convex optimization problem with theoretical convergence and global optimality guarantee; 3) Low-redundancy: thanks to the orthogonality objective, SOR is designed specifically to select less redundant features without sacrificing quality; 4) Extendability: SOR can enhance an existing set of preselected features by adding additional features that complement the existing feature set but still with strong predictive power. We present evaluation results showing that SOR consistently outperforms state of the art feature selection methods in a range of quality metrics on several real world data sets. We demonstrate a case study of a large-scale clinical application for predicting early onset of Heart Failure (HF) using real Electronic Health Records (EHRs) data of over 10K patients for over 7 years. Leveraging SOR, we are able to construct accurate and robust predictive models

and derive potential clinical insights.

Dijun Luo

The University of Texas at Arlington  
dijun.luo@gmail.com

Fei Wang, Jimeng Sun, Marianthi Marka  
IBM T.J. Watson Research Center  
fwang@us.ibm.com, jimeng@us.ibm.com,  
mmarkat@us.ibm.com

Jianying Hu, Shahram Ebadollahi  
IBM  
jyhu@us.ibm.com, ebad@us.ibm.com

#### CP11

##### **IntruMine: Mining Intruders in Untrustworthy Data of Cyber-Physical Systems**

A Cyber-Physical System (CPS) integrates physical (i.e., sensor) devices with cyber (i.e., informational) components to form a situation-aware system that responds intelligently to dynamic changes in real-world. It has wide application to scenarios of traffic control, environment monitoring and battlefield surveillance. This study investigates the specific problem of intruder mining in CPS: With a large number of sensors deployed in a designated area, the task is real time detection of intruders who enter the area, based on untrustworthy data. We propose a method called IntruMine to detect and verify the intruders. IntruMine constructs monitoring graphs to model the relationships between sensors and possible intruders, and computes the position and energy of each intruder with the link information from these monitoring graphs. Finally, a confidence rating is calculated for each potential detection, reducing false positives in the results. IntruMine is a generalized approach. Two classical methods of intruder detection can be seen as special cases of IntruMine under certain conditions. We conduct extensive experiments to evaluate the performance of IntruMine on both synthetic and real datasets and the experimental results show that IntruMine has better effectiveness and efficiency than existing methods.

Lu-An Tang, Quanquan Gu, Xiao Yu, Jiawei Han

UIUC  
tangl8@uiuc.edu, qgu3@illinois.edu, xiaoyu1@illinois.edu,  
hanj@illinois.edu

Thomas La Porta  
PSU  
tlp@cse.psu.edu

Alice Leung  
BBN Technology  
aleung@bbn.com

Tarek Abdelzaher  
UIUC  
zaher@illinois.edu

Lance Kaplan  
U.S. Army Lab  
lance.m.kaplan.civ@mail.mil

#### CP11

##### **Robust Reputation-Based Ranking on Bipartite**

## Rating Networks

With the growth of the Internet and E-commerce, bipartite rating networks are ubiquitous. In such bipartite rating networks, there exist two types of entities: the users and the objects, where users give ratings to objects. A fundamental problem in such networks is how to rank the objects by user's ratings. Although it has been extensively studied in the past decade, the existing algorithms either cannot guarantee convergence, or are not robust to the spammers. In this paper, we propose six new reputation-based algorithms, where the users' reputation is determined by the aggregated difference between the users' ratings and the corresponding objects' rankings. We prove that all of our algorithms converge into a unique fixed point. The time and space complexity of our algorithms are linear w.r.t. the size of the graph, thus they can be scalable to large datasets. Moreover, our algorithms are robust to the spamming users. We evaluate our algorithms using three real datasets. The experimental results confirm the effectiveness, efficiency, and robustness of our algorithms.

Rong-Hua Li, Jeffery Xu Yu, Xin Huang, Hong Cheng  
The Chinese University of Hong Kong  
rhli@se.cuhk.edu.hk, yu@se.cuhk.edu.hk,  
xhuang@se.cuhk.edu.hk, hcheng@se.cuhk.edu.hk

## CP11

### Mining Massive Archives of Mice Sounds with Symbolized Representations

The house mouse has long been an important model organism in biology and medicine to address human diseases. Advances in sensor technology have created a situation where our ability to collect data far outstrips our ability to analyze it manually. In this work we show a novel technique for mining mice vocalizations directly in the visual (spectrogram) space and the use of similarity search, classification, motif discovery and contrast set mining in this domain.

Jesin Zakaria  
3337 Utah Street  
Riverside, CA-92507  
jzaka001@ucr.edu

Sarah Rotschafer  
Department of Psychology  
University of California Riverside  
srots001@ucr.edu

Abdullah Mueen  
University of California, Riverside  
mueen@cs.ucr.edu

Khaleel Razak  
Department of Psychology  
University of California Riverside  
khaleel@ucr.edu

Eamonn Keogh  
University of California, Riverside  
eamonn@cs.ucr.edu

## MS1

### Efficient Monte Carlo Computation of Fisher In-

### formation Matrix using Prior Information

Abstract not available at time of publication.

Sonjoy Das  
University at Buffalo  
sonjoy@buffalo.edu

James Spall  
Johns Hopkins University  
james.spall@jhuapl.edu

Roger Ghanem  
University of Southern California  
Aerospace and Mechanical Engineering and Civil  
Engineering  
ghanem@usc.edu

## MS1

### Probabilistic Models of Past Climate Change

Abstract not available at time of publication.

Julien Emile-Geay, Dominique Guillot  
University of Southern California  
Los Angeles, CA 90089 0740, USA  
julieneg@usc.edu, dguillot@usc.edu

Tapio Schneider  
California Institute of Technology  
tapio@caltech.edu

Bala Rajaratnam  
Stanford University  
brajarat@stanford.edu

## MS1

### Diffusion on Random Manifolds

Abstract not available at time of publication.

Hadi Meidani  
University of Southern California  
meidani@usc.edu

Roger Ghanem  
University of Southern California  
Aerospace and Mechanical Engineering and Civil  
Engineering  
ghanem@usc.edu

## MS1

### A Priori Testing of Adaptive Sampling and Sparse PC Representations for Ocean General Circulation Models

Abstract not available at time of publication.

Justin Winokur  
Johns Hopkins University  
jwinokur@jhu.edu

Patrick R. Conrad  
MIT  
prconrad@mit.edu

Ihab Sraj, Alen Alexanderian  
Johns Hopkins University

israj@jhu.edu, aalex20@jhu.edu

Mohamed Iskandarani  
Rosenstiel School of Marine and Atmospheric Sciences  
University of Miami  
MIskandarani@rsmas.miami.edu

Ashwanth Srinivasan  
University of Miami  
asrinivasan@rsmas.miami.edu

Youssef M. Marzouk  
Massachusetts Institute of Technology  
ymarz@mit.edu

Omar Knio  
Duke University  
amk@duke.edu

### PP1

#### On Influential Node Discovery in Dynamic Social Networks

The problem of maximizing influence spread has been widely studied in social networks, because of its tremendous number of applications in determining critical points in a social network for information dissemination. All the techniques proposed in the literature are inherently static in nature, which are designed for social networks with a fixed set of links. However, many forms of *social interactions* are *transient* in nature, with relatively short periods of interaction. Any influence spread may happen only during the period of interaction, and the probability of spread is a function of the corresponding interaction time. Furthermore, such interactions are quite fluid and evolving, as a result of which the topology of the underlying network may change rapidly, as new interactions form and others terminate. In such cases, it may be desirable to determine the influential nodes based on the dynamic interaction patterns. Alternatively, one may wish to discover the most likely starting points for a *given infection pattern*. We will propose methods which can be used both for optimization of information spread, as well as the backward tracing of the source of influence spread. We will present experimental results illustrating the effectiveness of our approach on a number of real data sets.

Charu C. Aggarwal  
IBM T. J. Watson Research Center  
charu@us.ibm.com

Shuyang Lin, Philip Yu  
University of Illinois at Chicago  
slin38@cs.uic.edu, psyu@cs.uic.edu

### PP1

#### Event Detection in Social Streams

Social networks generate a large amount of text content over time because of continuous interaction between participants. The mining of such *social streams* is more challenging than traditional text streams, because of the presence of both text content and implicit network structure within the stream. The problem of event detection is also closely related to clustering, because the events can only be inferred from *aggregate* trend changes in the stream. In this paper, we will study the two related problems of clustering and event detection in social streams. We will study

both the supervised and unsupervised case for the event detection problem. We present experimental results illustrating the effectiveness of incorporating network structure in event discovery over purely content-based methods.

Charu C. Aggarwal  
IBM T. J. Watson Research Center  
charu@us.ibm.com

Karthik Subbian  
University of Minnesota  
karthik@umn.edu

### PP1

#### Query-based Biclustering using Formal Concept Analysis

Abstract not available at time of publication.

Faris Alqadah  
Johns Hopkins University  
faris.alqadah@gmail.com

Joel S. Bader  
Johns Hopkins University  
Department of Biomedical Engineering  
joel.bader@jhu.edu

Rajul Anand, Chandan Reddy  
Wayne State University  
rajulanand@wayne.edu, reddy@cs.wayne.edu

### PP1

#### Granger Causality Analysis in Irregular Time Series

In this paper, we propose a nonparametric generalization of the Granger graphical models called Generalized Lasso Granger (GLG) to uncover the temporal dependencies from *Irregular Time Series*, whose observations are not sampled at equally-spaced time stamps. Via theoretical analysis and extensive experiments, we verify the effectiveness of our model. Furthermore, we apply GLG to the application dataset of  $\delta^{18}O$  isotope of Oxygen records in Asia to discover the moisture transportation patterns in a 800-year period.

Mohammad Taha Bahadori, Yan Liu  
University of Southern California  
mohammab@usc.edu, yanliu.cs@usc.edu

### PP1

#### Clustering Based on Yukawa Potential

Inspired by the clustering phenomenon of nucleus, we propose a novel dynamic clustering algorithm based on Yukawa potential (Yupc). Each data object is regarded as a particle following the basic rules of movements in the Yukawa potential field. After several time intervals, similar objects gradually aggregate together and form clear clusters. Natural clusters of different shapes, densities, sizes, numbers and distributions can be detected by Yupc, reflecting the intrinsic structure of the original data set.

Xue Bai, Zezhen Lin, Yun Xiong, Yangyong Zhu  
Fudan University  
xuebai@fudan.edu.cn, justinlin722@gmail.com,  
yunx@fudan.edu.cn, yyzhu@fudan.edu.cn

PP1

**Balancing Prediction and Recommendation Accuracy: Hierarchical Latent Factors for Preference Data**

Recent works in Recommender Systems (RS) have investigated the relationships between the prediction accuracy, i.e. the ability of a RS to minimize cost functions in estimating users' preferences, and the accuracy of the recommendation list provided to users. Algorithms, which focus on the minimization of cost functions, have shown to achieve a weak recommendation accuracy, and vice versa. We present a Bayesian probabilistic hierarchical approach for RS designed to meet both prediction and recommendation accuracy.

Ettore Ritacco, [Nicola Barbieri](#), Giuseppe Manco, Riccardo Ortale  
ICAR-CNR  
ritacco@icar.cnr.it, barbieri@icar.cnr.it,  
manco@icar.cnr.it, ortale@icar.cnr.it

PP1

**Deterministic Cur for Improved Large-Scale Data Analysis: An Empirical Study**

Low-rank approximations which are computed from selected rows and columns of a given data matrix have attracted considerable attention lately. They have been proposed as an alternative to the SVD because they naturally lead to interpretable decompositions which was shown to be successful in application such as fraud detection, fMRI segmentation, and collaborative filtering. The CUR decomposition of large matrices, for example, samples rows and columns according to a probability distribution that depends on the Euclidean norm of rows or columns or on other measures of statistical leverage. At the same time, there are various deterministic approaches that do not resort to sampling and were found to often yield factorization of superior quality with respect to reconstruction accuracy. However, these are hardly applicable to large matrices as they typically suffer from high computational costs. Consequently, many practitioners in the field of data mining have abandoned deterministic approaches in favor of randomized ones when dealing with today's large-scale data sets. In this paper, we empirically disprove this prejudice. We do so by introducing a novel, linear-time, deterministic CUR approach that adopts the recently introduced Simplex Volume Maximization approach for column selection. The latter has already been proven to be successful for NMF-like decompositions of matrices of billions of entries. Our exhaustive empirical study on more than 30 synthetic and real-world data sets demonstrates that it is also beneficial for CUR-like decompositions. Compared to other deterministic CUR-like methods, it provides comparable reconstruction quality but operates much faster so that it easily scales to matrices of billions of elements. Compared to sampling-based methods, it provides competitive reconstruction quality while staying in the same run-time complexity class.

Christian Thureau, Kristian Kersting,  
[Christian Bauckhage](#)  
Fraunhofer IAIS  
cthureau@gmail.com, kristian.kersting@iais.fraunhofer.de,  
christian.bauckhage@iais.fraunhofer.de

PP1

**Combining Active Learning and Dynamic Dimen-****sionality Reduction**

To date, many active learning techniques have been developed for acquiring labels when training data is limited. However, an important aspect of the problem has often been neglected or just mentioned in passing: the curse of dimensionality. Yet, the curse of dimensionality poses even greater challenges in the case of limited data, which is precisely the setup for active learning. Reducing the dimensions is not a trivial task, however, as the correct number of dimensions depends on a number of factors including the training data size, the number of classes, the discriminative power of the features, and the underlying classification model. Moreover, active learning is typically applied in an iterative manner where the number of labels is smaller in the earlier iterations compared to the later ones. We propose an adaptive dimensionality reduction technique that determines the appropriate number of dimensions for each active learning iteration, utilizing the labeled and unlabeled data effectively to learn more accurate models. Extensive experiments comparing various approaches and parameter settings show that the proposed method improves performance drastically on three real-world text classification tasks.

[Mustafa Bilgic](#)  
Illinois Institute of Technology  
mbilgic@iit.edu

PP1

**Context-Aware Search for Personal Information Management Systems**

We present a novel context-aware desktop search framework by leveraging Hidden Markov Model to capture the relationships between user's access actions and activity states. The model is learned from user's past access history and is used to predict user's current activity upon the submission of some keyword query. We further propose a ranking scheme with this predicted context information incorporated. Experimental evaluation demonstrates its enhancement to user's search experience.

[Jidong Chen](#)  
EMC Research China  
EMC Corporation  
jidong.chen@emc.com

Wentao Wu  
Fudan University  
wentaowu@fudan.edu.cn

Hang Guo  
EMC Research China  
hang.guo@emc.com

Wei Wang  
Fudan University  
weiwang1@fudan.edu.cn

PP1

**Mining Social Dependencies in Dynamic Interaction Networks**

User-to-user interactions have become ubiquitous in Web 2.0. Users exchange emails, post on newsgroups, tag web pages, co-author papers, etc. Through these interactions, users co-produce or co-adopt content items (e.g., words in emails, tags in social bookmarking sites). We model such



dynamic interactions as a user interaction network, which relates users, interactions, and content items over time. After some interactions, a user may produce content that is more similar to those produced by other users previously. We term this effect *social dependency*, and we seek to mine from such networks the degree to which a user may be socially dependent on another user over time. We propose a *Decay Topic Model* to model the evolution of a user's preferences for content items at the topic level, as well as a *Social Dependency Metric* that quantifies the extent of social dependency based on interactions and content changes. Our experiments on two user interaction networks induced from real-life datasets show the effectiveness of our approach.

Freddy Chua, Hady Lauw, Ee-Peng Lim  
Singapore Management University  
freddycct@gmail.com, hadywlaww@smu.edu.sg,  
eplim@smu.edu.sg

### PP1

#### Detecting Irregularly Shaped Significant Spatial and Spatio-Temporal Clusters

Detecting significant overdensity or underdensity clusters in spatio-temporal data is critical for many real-world applications. Most existing approaches are designed to deal with regularly shaped clusters such as circular, elliptic and rectangular ones, but cannot work well on irregularly shaped clusters. In this paper, we propose GridScan, a grid-based approach for detecting irregularly shaped spatial clusters. In GridScan, a cluster is asymptotically described by a set of connected grid cells and is computed by a fast greedy region-growing algorithm with elaborating cluster merging in the process. The time complexity of GridScan is linear to the number of grids, making it scalable to very large datasets. A prospective spatio-temporal cluster detection approach, GridScan-Pro, is also proposed by extending GridScan. Experiments and a case study in the epidemic scenario demonstrate that our approaches greatly outperform existing ones in terms of accuracy, efficiency, and scalability.

Weishan Dong, Xin Zhang, Li Li, Changhua Sun, Lei Shi, Wei Sun  
IBM Research - China  
dongweis@cn.ibm.com, zxin@cn.ibm.com,  
lilichina@cn.ibm.com, schangh@cn.ibm.com,  
shllsh@cn.ibm.com, weisun@cn.ibm.com

### PP1

#### Contextual Collaborative Filtering Via Hierarchical Matrix Factorization

Matrix factorization (MF) has been demonstrated to be one of the most competitive techniques for collaborative filtering. However, state-of-the-art MFs do not consider contextual information, where ratings can be generated under different environments. For example, users select items under various situations, such as happy mood vs. sad, mobile vs. stationary, movies vs. book, etc. Under different contexts, the preference of users are inherently different. The problem is that MF methods uniformly decompose the rating matrix, and thus they are unable to factorize for different contexts. To amend this problem and improve recommendation accuracy, we introduce a "hierarchical" factorization model by considering the local context when performing matrix factorization. The intuition is that: as ratings are being generated from heterogeneous

environments, certain user and item pairs tend to be more similar to each other than others, and hence they ought to receive more collaborative information from each other. To take the contextual information into consideration, the proposed "contextual collaborative filtering" approach splits the rating matrix hierarchically by grouping similar users and items together, and factorizes each sub-matrix locally under different contexts. By building an ensemble model, the approach further avoids over-fitting with less parameter tuning. We analyze and demonstrate that the proposed method is a model-averaging gradient boosting model, and its error rate can be bounded. Experimental results show that it outperforms three state-of-the-art algorithms on a number of real-world datasets (MovieLens, Netflix, etc). The source code and datasets are available for download <http://www.cse.ust.hk/~ezhong/code/sdm12hmf.zip>.

Erheng Zhong  
Sun Yat-Sen University  
ezhong@cse.ust.hk

### Wei Fan

IBM T.J.Watson Research,  
weifan@us.ibm.com

Qiang Yang  
Department of Computer Science,  
Hong Kong University of Science  
qyang@cse.ust.hk

### PP1

#### Active Learning with Monotonicity Constraints

In many applications of data mining it is known beforehand that the response variable should be increasing (or decreasing) in the attributes. We propose two algorithms to exploit such monotonicity constraints for active learning in ordinal classification in two different settings. The basis of our approach is the observation that if the class label of an object is given, then the monotonicity constraints may allow the labels of other objects to be inferred. For instance, from knowing that loan applicant  $a$  is rejected, it can be concluded that all applicants that score worse than  $a$  on all criteria should be rejected as well. We propose two heuristics to determine good query points. These heuristics make a selection based on a point's potential to infer the labels of other points. The algorithms, each implemented with the proposed heuristics, are evaluated on artificial and real data sets to study their performance. We conclude that exploitation of monotonicity constraints can be very beneficial in active learning.

Ad Feelders, Nicola Barile  
Universiteit Utrecht  
A.J.Feelders@uu.nl, n.barile@uu.nl

### PP1

#### Pseudo Cold Start Link Prediction with Multiple Sources in Social Networks

Link prediction is an important task in social networks and data mining. Most existing researches therefore approach this problem by exploring the topological structure of the social network using only one source of information. In this work, we introduce the pseudo cold start link prediction with multiple source. We propose a two-phase supervised method. We assess our method empirically over a large

data collection obtained from Youtube.

Liang Ge

The State University of New York at Buffalo  
liange@buffalo.edu

### PP1

#### Discovering Context-Aware Influential Objects

It is very helpful for a user to get a *moderate* amount of information highly related to his/her immediate *context* (e.g., location, time, discussion topics) during the exploration of digital object collections (e.g., articles, web pages, blogs). For instance, in investigating a research topic, a researcher may be very interested in finding articles that are most related to the articles he/she already read on this topic, which we consider as ‘context’ in this paper. To facilitate users’ exploration, we introduce the problem of discovering Context-aware Influential Objects (CIO) from a collection of digital objects with influence relationships. Although there is a large amount of work in detecting direct influence degree between objects to denote how strong an object influences others, very few works utilize such direct influence to find influential objects for a context. To discover CIOs for a context consisting of several objects of a user’s interest, the first challenge is to meaningfully measure the *collective* influence of an object over a context considering both the direct influence and the *indirectly* derived influence, which is not taken into consideration by most ‘query by example’ approaches. We propose an aggregation framework to formulate the collective influence among objects by leveraging both direct and indirect influence. The second challenge is to discover CIOs *efficiently*. We present three approaches to calculate collective influence of an object over a context from an influence graph. In particular, the first approach utilizes the breadth-first-search paradigm; the other approaches make use of the topological sorting of graph nodes and perform context-aware search using *push* and *pull* mechanisms. We show experimental results on real datasets to demonstrate the effectiveness and efficiency of the proposed methodologies.

Huiping Cao

Computer Science  
New Mexico State University  
hcao@cs.nmsu.edu

Yangpai Liu, Yifan Hao

New Mexico State University  
lypmoon@nmsu.edu, yifan@nmsu.edu

Peng Han, Xinda Zeng

Chongqing Academy of Science and Technology, China  
han.peng@ciat-cq.org, shinda1020@gmail.com

### PP1

#### Monitoring and Mining Insect Sounds in Visual Space

Monitoring animals by the sounds they produce is an important and challenging task, whether the application is outdoors in a natural habitat, or in the controlled environment of a laboratory setting. In the former case the density and diversity of animal sounds can act as a measure of biodiversity. In the latter case, researchers often create control and treatment groups of animals, expose them to different interventions, and test for different outcomes. One possible manifestation of different outcomes may be changes in

the bioacoustics of the animals. With such a plethora of important applications, there have been significant efforts to build bioacoustic classification tools. However, we argue that most current tools are severely limited. They often require the careful tuning of many parameters (and thus huge amounts of training data), they are too computationally expensive for deployment in resource-limited sensors, they are specialized for a very small group of species, or they are simply not accurate enough to be useful. In this work we introduce a novel bioacoustic recognition/classification framework that mitigates or solves all of the above problems. We propose to classify animal sounds in the visual space, by treating the texture of their spectrograms as an acoustic fingerprint using a recently introduced parameter-free texture measure as a distance measure. We further show that by searching for the most representative acoustic fingerprint we can significantly outperform other techniques in terms of speed and accuracy.

Yuan Hao, Bilson J. Campana, Eamonn Keogh

University of California, Riverside

yhao@cs.ucr.edu,

bcampana@cs.ucr.edu,

eamonn@cs.ucr.edu

### PP1

#### Image Mining of Historical Manuscripts to Establish Provenance

The recent digitization of more than twenty million books has been led by initiatives from countries wishing to preserve their cultural heritage and by commercial endeavors, such as the Google Print Library Project. Within a few years a significant fraction of the world’s books will be online. For millions of intact books and tens of millions of loose pages, the provenance of the manuscripts may be in doubt or completely unknown, thus denying historians an understanding of the context of the content. In some cases it may be possible for human experts to regain the provenance by examining linguistic, cultural and/or stylistic clues. However, such experts are rare and this investigation is clearly a time-consuming process. One technique used by experts to establish provenance is the examination of the ornate initial letters appearing in the questioned manuscript. By comparing the initial letters in the manuscript to annotated initial letters whose origin is known, the provenance can be determined. In this work we show for the first time that we can reproduce this ability with a computer algorithm. We leverage off a recently introduced technique to measure texture similarity and show that it can recognize initial letters with an accuracy that rivals or exceeds human performance. A brute force implementation of this measure would require several years to process a single large book; however, we introduce a novel lower bound that allows us to process the books in minutes.

Bing Hu

University of California, Riverside

University of California, Riverside

bhu002@ucr.edu

### PP1

#### RP-growth: Top-k Mining of Relevant Patterns with Minimum Support Raising

This paper proposes RP-growth, an efficient top-k mining algorithm for the patterns highly relevant to the class of interest. RP-growth conducts branch-and-bound search using anti-monotonic upper bounds of the relevance score,

and its pruning strategy is successfully translated to minimum support raising, a standard pruning strategy in top-k mining. Furthermore, RP-growth introduces an aggressive pruning strategy based on the notion called weakness. Experimental results on text classification exhibit the efficiency and the usefulness of RP-growth.

Yoshitaka Kameya, Taisuke Sato  
Tokyo Institute of Technology  
kameya@mi.cs.titech.ac.jp, sato@mi.cs.titech.ac.jp

## PP1

### Fast Random Walk Graph Kernel

Random walk graph kernel has been used as an important tool for various data mining tasks including classification and similarity computation. Despite its usefulness, however, it suffers from the expensive computational cost which is at least  $O(n^3)$  or  $O(m^2)$  for graphs with  $n$  nodes and  $m$  edges. In this paper, we propose Ark, a set of fast algorithms for random walk graph kernel computation. Ark is based on the observation that real graphs have much lower intrinsic ranks, compared with the orders of the graphs. Ark exploits the low rank structure to quickly compute random walk graph kernels in  $O(n^2)$  or  $O(m)$  time. Experimental results show that our method is up to 97,865 times faster than the existing algorithms, while providing more than 91.3% of the accuracies.

U Kang  
Carnegie Mellon University  
Computer Science Department  
ukang@cs.dot.cmu.dot.edu

Hanghang Tong  
IBM T.J. Watson  
htong@us.ibm.com

Jimeng Sun  
IBM T.J. Watson Research Center  
jimeng@us.ibm.com

## PP1

### Tracking Spatio-Temporal Diffusion in Climate Data

A forest canopy forms a critical platform for complex interactions between the vegetation and the atmosphere boundary layer and is considered as a crucial piece for environmental scientists in their understanding of the ecosystem and its response to the climate change. Microfronts represent a class of these interactions characterized by a moving mass of air that introduce fluctuations in ambient temperature and humidity on small spatial and temporal scales. In this paper, we present a joint spatio-temporal hidden markov model that simultaneously incorporates neighborhood dependencies in space and time. We show that our approach can trace the diffusion of microfronts more effectively than several baseline methods over a sensor data from Brazilian rainforest and a synthetically generated dataset.

Jaya Kawale, Aditya Pal  
University of Minnesota  
kawale@cs.umn.edu, apal@cs.umn.edu

Rob Fatland  
Microsoft Research  
rob.fatland@microsoft.com

## PP1

### Group Sparsity in Nonnegative Matrix Factorization

A recent challenge in data analysis for science and engineering is that data are often represented in a structured way. In particular, many data mining tasks have to deal with group-structured prior information, where features or data items are organized into groups. In this paper, we develop group sparsity regularization methods for nonnegative matrix factorization (NMF). NMF is an effective data mining tool that has been widely adopted in text mining, bioinformatics, and clustering, but a principled approach to incorporating group information into NMF has been lacking in the literature. Motivated by an observation that features or data items within a group are expected to share the same sparsity pattern in their latent factor representation, we propose mixed-norm regularization to promote group sparsity in the factor matrices of NMF. Group sparsity improves the interpretation of latent factors. Efficient convex optimization methods for dealing with the mixed-norm term are presented along with computational comparisons between them. Application examples of the proposed method in factor recovery, semi-supervised clustering, and multilingual text analysis are demonstrated.

Jingu Kim  
Georgia Institute of Technology  
jingu@cc.gatech.edu

Renato C. Monteiro  
Georgia Institute of Technology  
School of ISyE  
monteiro@isye.gatech.edu

Haesun Park  
Georgia Institute of Technology  
hpark@cc.gatech.edu

## PP1

### Global Linear Neighborhoods for Efficient Label Propagation

Graph-based semi-supervised learning improves classification by combining labeled and unlabeled data through label propagation. In this paper, we propose to learn a non-negative low-rank graph to capture global linear neighborhoods, under the assumption that each data point can be linearly reconstructed from weighted combinations of its direct neighbors and reachable indirect neighbors. Large scale experiments on UCI datasets and gene expression datasets showed label propagation based on global linear neighborhoods achieved more accurate classification results.

Ze Tian, Rui Kuang  
Dept Computer Science  
University of Minnesota  
tianze@cs.umn.edu, kuang@cs.umn.edu

## PP1

### Generalized Similarity Kernels for Efficient Sequence Classification

String kernel-based machine learning methods have yielded great success in practical tasks of structured/sequential data analysis such as document topic elucidation, music genre classification, protein superfamily and fold prediction. However, typical string kernel methods rely on *sym-*

*bold Hamming-distance* based matching which may not necessarily reflect the underlying (e.g., physical) similarity between sequence fragments. In this work we propose a novel computational framework that uses more “precise”, *general similarity metrics*  $\mathcal{S}(\cdot, \cdot)$  and distance-preserving embeddings with string kernels and improves upon state-of-the-art on a number of sequence analysis tasks such as music, and biological sequence classification.

Pavel P. Kuksa  
Rutgers University  
pkuksa@nec-labs.com

Imdadullah Khan  
Gulf University for Science and Technology  
imdadk@gmail.com

Vladimir Pavlovic  
Rutgers University  
vladimir@cs.rutgers.edu

### PP1

#### Detecting Extreme Rank Anomalous Collections

Anomaly or outlier detection has a wide range of applications, including fraud and spam detection. Most existing studies focus on detecting point anomalies, i.e., individual, isolated entities. However, there is an increasing number of applications in which anomalies do not occur individually, but in small collections. Unlike the majority, entities in an anomalous collection tend to share certain extreme behavioral traits. The knowledge essential in understanding why and how the set of entities becomes outliers would only be revealed by examining at the collection level. A good example is web spammers adopting common spamming techniques. To discover this kind of anomalous collections, we introduce a novel definition of anomaly, called *Extreme Rank Anomalous Collection*. We propose a statistical model to quantify the anomalousness of such a collection, and present an exact as well as a heuristic algorithms for finding top- $K$  extreme rank anomalous collections. We apply the algorithms on real Web spam data to detect spamming sites, and on IMDB data to detect unusual actor groups. Our algorithms achieve higher precisions compared to existing spam and anomaly detection methods. More importantly, our approach succeeds in finding meaningful anomalous collections in both datasets.

Hanbo Dai, Feida Zhu, Ee-Peng Lim, Hwee Hwa Pang,  
Hady Lauw  
Singapore Management University  
hanbo.dai.2008@smu.edu.sg, fdzhu@smu.edu.sg,  
eplim@smu.edu.sg, hhpang@smu.edu.sg,  
hadywlaw@smu.edu.sg

### PP1

#### Visualizing Variable-Length Time Series Motifs

The problem of time series motif discovery has received a lot of attention from researchers in the past decade. Most existing work on finding time series motifs require that the length of the motifs be known in advance. However, such information is not always available. In addition, motifs of different lengths may co-exist in a time series dataset. In this work, we develop a motif visualization system based on grammar induction. We demonstrate that grammar induction in time series can effectively identify repeated patterns without prior knowledge of their lengths. The motifs dis-

covered by the visualization system are of variable lengths in two ways. Not only can the *inter-motif* subsequences be of different lengths, the *intra-motif* subsequences also are not restricted to have identical length—a unique property that is desirable, but has not been seen in the literature.

Yuan Li, Jessica Lin  
George Mason University  
Department of Computer Science  
ylif@gmu.edu, jessica@cs.gmu.edu

Tim Oates  
Department of Computer Science and Electrical  
Engineering  
University of Maryland Baltimore County, Baltimore,  
USA  
oates@cs.umbc.edu

### PP1

#### Which Distance Metric Is Right: An Evolutionary K-Means View

We study the impact of monotone metrics on K-means clustering. By revealing the order-preserving property and proving that the cluster centroid is a good approximator of their respective optimal centers, we show K-means cannot differentiate the cosine-monotone metrics. Then an evolutionary framework is proposed to enable inspection of these metrics. Most importantly, this paper furthers our understanding of the impact of the metrics on the optimization process of K-means.

Chuanren Liu  
Rutgers, the State University of New Jersey  
chuanren.liu@rutgers.edu

Tianming Hu  
Dongguan University of Technology  
tmhu@ieee.org

Yong Ge, Hui Xiong  
Rutgers, the State University of New Jersey  
yongge@rutgers.edu, hxiong@rutgers.edu

### PP1

#### Constructing Training Sets for Outlier Detection

Outlier detection often works in an unsupervised manner due to the difficulty of obtaining enough training data. Since outliers are rare, one has to label a very large dataset to include enough outliers in the training set, with which classifiers could sufficiently learn the concept of outliers. Labeling a large training set is costly for most applications. However, we could just label suspected instances identified by unsupervised methods. In this way, the number of instances to be labeled could be greatly reduced. Based on this idea, we propose CISO, an algorithm Constructing training set by Identifying Suspected Outliers. In this algorithm, instances in a pool are first ranked by an unsupervised outlier detection algorithm. Then, suspected instances are selected and hand-labeled, and all remaining instances receive label of inlier. As such, all instances in the pool are labeled and used in the training set. We also propose Budgeted CISO (BCISO), with which user could set a fixed budget for labeling. Experiments show that both algorithms achieve good performance compared to other methods when the same amount of labeling effort

are used.

Liping Liu  
EECS Oregon State University  
liping.liulp@gmail.com

Xiaoli Z. Fern  
Oregon State University  
xfern@eecs.oregonstate.edu

#### PP1

##### **A Flexible Open-Source Toolbox for Scalable Complex Graph Analysis**

The Knowledge Discovery Toolbox (KDT) enables domain experts to perform complex analyses of huge datasets on supercomputers using a high-level language without grappling with the difficulties of writing parallel code, calling parallel libraries, or becoming a graph expert. KDT provides a flexible Python interface to a small set of high-level reusable graph operations; composing these operations produces graph analysis algorithms scalable to graphs on the order of 10 billion edges or greater.

Adam Lugowski  
UC Santa Barbara  
alugowski@cs.ucsb.edu

Aydin Buluc  
Lawrence Berkeley National Laboratory  
abuluc@lbl.gov

David Alber  
Microsoft  
david.alber@microsoft.com

John R. Gilbert  
Dept of Computer Science  
University of California, Santa Barbara  
gilbert@cs.ucsb.edu

Steve Reinhardt  
Cray, Inc.  
spr@cray.com

Yun Teng, Andrew Waranis  
UC Santa Barbara  
yunteng@umail.ucsb.edu, andrewwaranis@umail.ucsb.edu

#### PP1

##### **Fast Robustness Estimation in Large Social Graphs: Communities and Anomaly Detection**

Given a large social graph, what can we say about its robustness? In this work, we are trying to answer the above question studying the *expansion properties* of large social graphs. We present a measure which characterizes the robustness of a graph and serves as global measure of the community structure (or lack thereof), and we show how to compute it efficiently. We present extensive experimental results on both static and time-evolving real networks.

Fragkiskos D. Malliaros  
Department of Computer Engineering and Informatics  
University of Patras  
malliaro@ceid.upatras.gr

Vasileios Megalooikonomou  
Department of Computer Engineering and Informatics  
University of Patras and Temple University  
vasilis@ceid.upatras.gr

Christos Faloutsos  
Carnegie Mellon University  
christos@cs.cmu.edu

#### PP1

##### **On Finding Joint Subspace Boolean Matrix Factorizations**

Abstract not available at time of publication.

Pauli Miettinen  
Max Planck Institute for Informatics  
pauli.miettinen@mpi-inf.mpg.de

#### PP1

##### **Generalized Optimization Framework for Graph-Based Semi-Supervised Learning**

We develop a generalized optimization framework for graph-based semi-supervised learning. The framework gives as particular cases the Standard Laplacian, Normalized Laplacian and PageRank based methods. We have also provided new probabilistic interpretation based on random walks and characterized the limiting behaviour of the methods. The random walk based interpretation allows us to explain differences between the performances of methods with different smoothing kernels. It appears that the PageRank based method is robust with respect to the choice of the regularization parameter and the labelled data. We illustrate our theoretical results with two realistic datasets, characterizing different challenges: Les Misérables characters social network and Wikipedia hyper-link graph. The graph-based semi-supervised learning classifies the Wikipedia articles with very good precision and perfect recall employing only the information about the hyper-text links.

Marina M. Sokol, Konstantin Avrachenkov  
INRIA Sophia Antipolis  
marina.sokol@inria.fr, k,avrachenkov@sophia.inria.fr

Paulo Goncalves  
INRIA Rhone  
paulo.goncalves@ens-lyon.fr

Alexey Mishenin  
St. Petersburg State University  
alexey.mishenin@gmail.com

#### PP1

##### **A Tree-Based Kernel for Graphs**

This paper proposes a new tree-based kernel for graphs. Graphs are decomposed into multisets of ordered Directed Acyclic Graphs (DAGs) and a family of kernels is defined by application of tree kernels extended to the DAG domain. We focus our attention on the efficient development of one member of this family. A technique for speeding up the computation is given, as well as theoretical bounds and practical evidence of its feasibility.

Nicolo' Navarin, Giovanni Da San Martino, Alessandro Sperduti

University of Padova  
 Department of Mathematics  
 nnavarin@math.unipd.it,  
 sperduti@math.unipd.it

dasan@math.unipd.it,

### PP1

#### Density-Based Projected Clustering over High Dimensional Data Streams

Clustering of high dimensional data streams is an important problem in many application domains, a prominent example being network monitoring. Several approaches have been lately proposed for solving independently the different aspects of the problem. There exist methods for clustering over full dimensional streams and methods for finding clusters in subspaces of high dimensional static data. Yet only a few approaches have been proposed so far which tackle both the stream and the high dimensionality aspects of the problem simultaneously. In this work, we propose a new density-based projected clustering algorithm, HDDStream, for high dimensional data streams. Our algorithm summarizes both the data points and the dimensions where these points are grouped together and maintains these summaries online, as new points arrive over time and old points expire due to ageing. Our experimental results illustrate the effectiveness and the efficiency of HDDStream and also demonstrate that it could serve as a trigger for detecting drastic changes in the underlying stream population, like bursts of network attacks.

Eirini C. Ntoutsi

Ludwig-Maximilians-Universität München (LMU)  
 ntoutsi@dbs.ifi.lmu.de

Arthur Zimek  
 LMU Munich  
 zimek@dbs.ifi.lmu.de

Themis Palpanas  
 Information Engineering and Computer Science  
 Department  
 (DISI), University of Trento, Italy  
 themis@disi.unitn.eu

Peer Kröger  
 Ludwig-Maximilians-Universität München  
 kroeger@dbs.ifi.lmu.de

Hans-Peter Kriegel  
 Ludwig-Maximilians University Munich  
 kriegel@dbs.ifi.lmu.de

### PP1

#### A Novel Approximation to Dynamic Time Warping Allows Anytime Clustering of Massive Time Series Datasets

Given the ubiquity of time series data, the data mining community has spent significant time investigating the best time series similarity measure to use for various tasks and domains. After more than a decade of extensive efforts, there is increasing evidence that Dynamic Time Warping (DTW) is very difficult to beat. Given that, recent efforts have focused on making the intrinsically slow DTW algorithm faster. For the similarity-search task, an important subroutine in many data mining algorithms, significant progress has been made by replacing the vast majority of expensive DTW calculations with cheap-to-compute lower

bound calculations. However, these lower bound based optimizations do not directly apply to clustering, and thus for some realistic problems, clustering with DTW can take days or weeks. In this work, we show that we can mitigate this untenable lethargy by casting DTW clustering as an anytime algorithm. At the heart of our algorithm is a novel data-adaptive approximation to DTW which can be quickly computed, and which produces approximations to DTW that are much better than the best currently known linear-time approximations. We demonstrate our ideas on real world problems showing that we can get virtually all the accuracy of a batch DTW clustering algorithm in a fraction of the time.

Qiang Zhu, Gustavo E. Batista,

Thanawin Rakthanmanon, Emaonn Keogh

University of California, Riverside

qzhu@cs.ucr.edu,

gbatista@cs.ucr.edu,

rakthant@cs.ucr.edu, eamonn@cs.ucr.edu

### PP1

#### Nearest-Neighbor Search on a Time Budget via Max-Margin Trees

Many high-profile applications pose high-dimensional nearest-neighbor search problems. Yet, it still remains difficult to achieve fast query times for state-of-the-art approaches which use multidimensional trees for either exact or approximate search, possibly in combination with hashing approaches. Moreover, a number of these applications only have a limited amount of time to answer nearest-neighbor queries. However, we observe empirically that the correct neighbor is often found early within the tree-search process, while the bulk of the time is spent on verifying its correctness. Motivated by this, we propose an algorithm for finding the best neighbor given any particular time limit, and develop a new data structure, the *max-margin tree*, to achieve accurate results even with small time budgets. Max-margin trees perform better in the limited-time setting than current commonly-used data structures such as the *kd-tree* and the more recently developed *RP-tree* data structure.

Parikshit Ram

School of Computational Science and Engineering

Georgia Institute of Technology

p.ram@gatech.edu

Dongryeol Lee, Alexander Gray

Georgia Institute of Technology

dongryel@cc.gatech.edu, agray@cc.gatech.edu

### PP1

#### Efficient Clustering of Metagenomic Sequences Using Locality Sensitive Hashing

The new generation of genomic technologies have allowed researchers to determine the collective DNA of organisms (e.g. microbes) co-existing as communities across the ecosystem (e.g. within the human host). There is a need for the computational approaches to analyze and annotate the large volumes of available sequence data from such microbial communities (metagenomes). In this paper, we developed an efficient and accurate metagenome clustering approach that uses the locality sensitive hashing (LSH) technique to approximate the computational complexity associated with comparing sequences. We introduce the use of fixed-length, gapless subsequences for improving the sensitivity of the LSH-based similarity func-

tion. We evaluate the performance of our algorithm on two metagenome datasets associated with microbes existing across different human skin locations. Our empirical results show the strength of the developed approach in comparison to three state-of-the-art sequence clustering algorithms with regards to computational efficiency and clustering quality. We also demonstrate practical significance for the developed clustering algorithm, to compare bacterial diversity and structure across different skin locations.

Zeehasham Rasheed, Huzefa Rangwala, Daniel Barbara  
George Mason University  
zrasheed@gmu.edu, rangwala@cs.gmu.edu,  
dbarbara@gmu.edu

### PP1

#### On Evaluation of Outlier Rankings and Outlier Scores

Outlier detection research is currently focusing on the development of new methods and on improving the computation time for these methods. Evaluation however is rather heuristic, often considering just precision in the top  $k$  results or using the area under the ROC curve. These evaluation procedures do not allow for assessment of similarity between methods. Judging the similarity of or correlation between two rankings of outlier scores is an important question in itself but it is also an essential step towards meaningfully building outlier detection ensembles, where this aspect has been completely ignored so far. In this study, our generalized view of evaluation methods allows both to evaluate the performance of existing methods as well as to compare different methods w.r.t. their detection performance. Our new evaluation framework takes into consideration the class imbalance problem and offers new insights on similarity and redundancy of existing outlier detection methods. As a result, the design of effective ensemble methods for outlier detection is considerably enhanced.

Arthur Zimek, Erich Schubert, Remigius Wojdanowski  
LMU Munich  
zimek@dbf.lmu.de, schube@dbf.lmu.de,  
wojdanowski@dbf.lmu.de

Hans-Peter Kriegel  
Ludwig-Maximilians University Munich  
kriegel@dbf.lmu.de

### PP1

#### Regularized Structured Output Learning with Partial Labels

We consider the problem of learning structured output probabilistic models with training examples having partial labels. Partial label scenarios arise commonly in web applications such as hierarchical and multi-label classification. We solve the learning problem with partial labels by incorporating entropy and label distribution or correlation regularizations along with marginal likelihood maximization. We develop probabilistic taxonomy and multi-label classifier models, and provide the ideas needed for expanding their usage to the partial labels scenario.

Sundararajan Sellamanickam, Charu Tiwari  
Yahoo! Labs, Bangalore, India  
ssrajan@yahoo-inc.com, charu@yahoo-inc.com

Sathiya Keerthi Selvaraj  
Yahoo! Labs, Santa Clara, CA  
selvarak@yahoo-inc.com

### PP1

#### The Similarity Between Stochastic Kronecker and Chung-Lu Graph Models

The *Stochastic Kronecker Graph* (SKG) model has been chosen as a benchmark by the Graph500 steering committee, but there is little understanding of the properties of this model. We show that the parallel variant of the edge-configuration model given by Chung and Lu (CL) is similar to the SKG model. Our experiments suggest that the graph distribution represented by SKG is almost the same as that given by a CL model. Also, CL appears to fit real data as well as SKG.

C. Seshadhri, Ali Pinar  
Sandia National Labs  
scomand@sandia.gov, apinar@sandia.gov

Tamara G. Kolda  
Sandia National Laboratories  
tgkolda@sandia.gov

### PP1

#### Wigm: Discovery of Subgraph Patterns in a Large Weighted Graph

Many research areas have begun representing massive data sets as very large graphs. Thus, graph mining has been an active research area in recent years. Most of the graph mining research focuses on mining unweighted graphs. However, weighted graphs are actually more common. The weight on an edge may represent the likelihood or logarithmic transformation of likelihood of the existence of the edge or the strength of an edge, which is common in many biological networks. In this paper, a weighted subgraph pattern model is proposed to capture the importance of a subgraph pattern and our aim is to find these patterns in a large weighted graph. Two related problems are studied in this paper: (1) discovering all patterns with respect to a given minimum weight threshold and (2) finding  $k$  patterns with the highest weights. The weighted subgraph patterns do not possess the anti-monotonic property and in turn, most of existing subgraph mining methods could not be directly applied. Fortunately, the **1-extension** property is identified so that a bounded search can be achieved. A novel weighted graph mining algorithm, namely WIGM, is devised based on the 1-extension property. Last but not least, real and synthetic data sets are used to show the effectiveness and efficiency of our proposed models and algorithms.

Wei Su, Jiong Yang  
Case Western Reserve University  
wei.su@case.edu, jiong.yang@case.edu

Shirong Li  
Aliyun Inc  
shirong.li@alibaba-inc.com

Mehmet Dalkilic  
Indiana University  
dalkilic@indiana.edu

**PP1****Legislative Prediction Via Random Walks over a Heterogeneous Graph**

In this article, we propose a random walk-based model to predict legislators votes on a set of bills. In particular, we first convert roll call data, i.e. the recorded votes and the corresponding deliberative bodies, to a heterogeneous graph, where both the legislators and bills are treated as vertices. Three types of weighted edges are then computed accordingly, representing legislators social and political relations, bills semantic similarity, and legislator-bill vote relations. Through performing two-stage random walks over this heterogeneous graph, we can estimate legislative votes on past and future bills. We apply this proposed method on real legislative roll call data of the United States Congress and compare to state-of-the-art approaches. The experimental results demonstrate the superior performance and unique prediction power of the proposed model.

Jun Wang

IBM Thomas J. Watson Research Center  
Business Analytics and Mathematical Sciences  
Department  
wangjun@us.ibm.com

Kush Varshne, Aleksandra Mojsilovic  
IBM Thomas J. Watson Research Center  
krvarshn@us.ibm.com, aleksand@us.ibm.com

**PP1****An Iterative and Re-Weighting Framework for Rejection and Uncertainty Resolution in Crowdsourcing**

We propose an Iterative Re-weighted Consensus Maximization framework to address the missing and uncertain label problem. The intuitive idea is to use an iterated framework to estimate each labeler's hidden competence and formulate it as a spectral clustering problem in the functional space, in order to minimize the overall loss given missing and uncertain information. One main advantage of the proposed method from state-of-the-art Bayesian model averaging based approaches is that it uncovers the intrinsic consistency among different set of answers and mines the best possible ground truth.

Sihong Xie

University of Illinois at Chicago  
sxie6@uic.edu

Wei Fan  
IBM T.J.Watson Research  
wei.fan@gmail.com

Philip Yu  
University of Illinois at Chicago  
psyu@cs.uic.edu

**PP1****Citation Prediction in Heterogeneous Bibliographic Networks**

To reveal information hiding in link space of bibliographical networks, link analysis has been studied from different perspectives in recent years. In this paper, we address a novel problem namely citation prediction, that is: given information about authors, topics, target publication venues as well as time of certain research paper, finding and predict-

ing the citation relationship between a query paper and a set of previous papers. Considering the gigantic size of relevant papers, the loosely connected citation network structure as well as the highly skewed citation relation distribution, citation prediction is more challenging than other link prediction problems which have been studied before. By building a meta-path based prediction model on a topic discriminative search space, we here propose a two-phase citation probability learning approach, in order to predict citation relationship effectively and efficiently. Experiments are performed on real-world dataset with comprehensive measurements, which demonstrate that our framework has substantial advantages over commonly used link prediction approaches in predicting citation relations in bibliographical networks.

Xiao Yu, Quanquan Gu

UIUC  
xiaoyu1@illinois.edu, qgu3@illinois.edu

Mianwei Zhou

University of Illinois at Urbana Champaign  
zhou18@illinois.edu

Jiawei Han

UIUC  
hanj@illinois.edu

**PP1****Mining Multi-Label Data Streams Using Ensemble-Based Active Learning**

Data stream classification has drawn increasing attention from the data mining community in recent years, where a large number of stream classification models were proposed. However, most existing models were merely focused on mining from single-label data streams. Mining from multi-label data streams has not been fully addressed yet. On the other hand, although some recent work touched the multi-label stream mining problem, they never consider the expensive labeling cost issue, preventing them from real-world applications. To this end, we study, in this paper, a challenging problem that mining from multi-label data streams with limited labeling resource. Specifically, we propose an ensemble-based active learning framework to handle the large volume of stream data, expensive labeling cost and concept drifting problems on multi-label data streams. Experiments on both synthetic and real world data sets demonstrate the performance of the proposed method.

Peng Wang, Peng Zhang, Li Guo  
Institute of Information Engineering  
Chinese Academy of Sciences  
peng860215@gmail.com, zhangpeng04@gmail.com,  
guoli@ict.ac.cn

**PP1****Feature Selection for High-Dimensional Integrated Data**

Motivated by the problem of identifying correlations between genes or features of two related biological systems, we propose a model of *feature selection* in which only a subset of the predictors  $X_i$  are dependent on the multidimensional variate  $Y$ , and the remainder of the predictors constitute a "noise set"  $X_u$  independent of  $Y$ . Using Monte Carlo simulations, we investigated the relative performance of two methods: thresholding and singular-value decompo-



sition, in combination with stochastic optimization to determine ‘empirical bounds’ on the small-sample accuracy of an asymptotic approximation. We demonstrate utility of the thresholding and SVD feature selection methods with respect to a recent infant intestinal gene expression and metagenomics dataset.

Charles Y. Zheng  
Texas A and M  
charles.y.zheng@gmail.com

Scott Schwartz  
Texas Agrilife  
sschwartz@ag.tamu.edu

Robert Chapkin  
Texas A and M  
Integrative Nutrition and Complex Diseases  
r-chapkin@tamu.edu

Raymond Carroll  
Texas A and M  
Statistics  
carroll@stat.tamu.edu

Ivan Ivanov  
Texas A and M  
Veterinary Physiology and Pharmacology  
iivanov@cvm.tamu.edu