10/13/2018

Michigan Institute for Data Science (MIDAS)

Computational Challenges & Research Opportunities

Ivo D Dinov

www.MIDAS.umich.edu

Michigan Institute for Data Science (MIDAS) University of Michigan

National Big Data Science Curricula Constellation





MIDAS Faculty Affiliates Network

200+ Faculty Affiliates (Ann Arbor + Flint + Dearborn campuses)

Bio/clinical Informatics





Visual Analytics



ss Analytics













2

MIDAS Student Organizations

Computational Social Science Rackham Interdisciplinary Workshop

- 160+ members from 17 depts / institutes
- Panels, skill building workshops, reading groups, roundtable discussions

sites.lsa.umich.edu/css

Statistics in the Community (STATCOM)

"promotes student-driven programs that provide statistical consulting as a community service"

- 75+ members from Biostatistics, Statistics and Survey Methodology
- Offers services to local governmental and nonprofit community groups

sph.umich.edu/biostat/statcom

Michigan Data Science Team (MDST)

"teaches practical data science skills by solving impactful problems"

- 50+ active members from CoE, LS&A, Ross and other units
- Competitions, Tutorials, Projects <u>midas.umich.edu/mdst</u>

Michigan Student Artificial Intelligence Lab (MSAIL)

- 50+ members from CoE, LS&A, Ross and other units
- Machine learning reading group, research projects, tutorials <u>http://msail.org</u>







Data Science and Predictive Analytics (HS650)

Areas	Competency	Expectation			
	Tools	Working knowledge of basic software tools (command-line, GUI based, or web-services)			
Algorithms and Applications	Algorithms	Knowledge of core principles of scientific computing, applications programming, API's, algorithm complexity, an data structures			
	Application Domain	Data analysis experience from at least one application area, either through coursework, internship, research project, etc.			
Data Manage- ment	Data validation & visualization	Curation, Exploratory Data Analysis (EDA) and visualizati			
	Data wrangling	Skills for data normalization, data cleaning, data aggregation, and data harmonization/registration			
	Data infrastructure	Handling databases, web-services, Hadoop, multi-source data			
Analysis Methods	Statistical inference	Basic understanding of bias and variance, principles of (non)parametric statistical inference, and (linear) modeling			
	Study design and diagnostics	Design of experiments, power calculations and sample sizing, strength of evidence, p-values, False Discovery Rates			
	Machine Learning	Dimensionality reduction, k-nearest neighbors, random forests, AdaBoost, kernelization, SVM, ensemble methods, CNN			

Big Data Science Challenges & Opportunities



MIDAS Challenge Initiatives

Data-Intensive Transportation Research Hub

- Reinventing Public Urban Transportation and Mobility
- Building a Transportation Data Ecosystem

Data-Intensive Learning Analytics Hub

- LEAP: Analytics for LEarners As People
- HOME: Holistic Modeling of Education

Data-Intensive Social Science Research Hub

- Computational Approaches for the Construction of Novel Macroeconomic Data
- A Social Science Collaboration for Research on Communication and Learning based upon Big Data

Data-Intensive Health Science Research Hub

- Michigan Center for Single-Cell Genomic Data Analytics
- Michigan Integrated Center for Health Analytics & Medical Prediction (MiCHAMP)
- Identifying Real-Time Data Predictors of Stress and Depression Using Mobile Technology

Data Science for Music Hub

- Understanding how the brain processes music through the Bach trio sonatas
- Mining patterns of audience engagement / crowdsourced music performances
- The sound of text
- A computational study of patterned melodic structures across musical cultures

http://midas.umich.edu/research



Characteristics of Big Data

IBM Big Data 4V's: Volume, Variety, Velocity & Veracity

Big Bio Data Dimensions	Tools	Biomed Exa observationa
Size	Harvesting and management of vast amounts of data	on 10,000's s
Complexity	Wranglers for dealing with heterogeneous data	genetics, clin
Incongruency	Tools for data harmonization and aggregation	elements
Multi-source	Transfer and joint modeling of disparate elements	Software dev training, serv
Multi-scale	Macro to meso to micro scale observations	methodologic associated w
Time	Techniques accounting for longitudinal effects	Discovery Sc existing oppo
Incomplete	Reliable management of missing data	educators, re practitioners

Biomed Example: analyzing observational data of 1,000's Parkinson's disease patients based on 10,000's signature biomarkers derived from multi-source imaging, genetics, clinical, physiologic, obhenomics and demographic data elements

Software developments, student training, service platforms and methodological advances associated with the Big Data Discovery Science all present existing opportunities for learners, educators, researchers, practitioners and policy makers

Dinov, GigaScience (2016) PMID:26918190







Case-Studies – Parkinson's Disease

- Investigate falls in PD patients using clinical, demographic and neuroimaging data from two independent initiatives (UMich & Tel Aviv U)
- Applied <u>controlled feature selection</u> to identify the most salient predictors of patient falls (gait speed, Hoehn and Yahr stage, postural instability and gait difficulty-related measurements)
- Model-based (e.g., GLM) and model-free (RF, SVM, Xgboost) analytical methods used to forecasts clinical outcomes (e.g., falls)
- □ Internal statistical cross <u>validation</u> + external out-of-bag validation
- Four specific <u>challenges</u>
 - Challenge 1, harmonize & aggregate complex, multisource, multisite PD data
 - □ Challenge 2, identify salient predictive features associated with specific clinical traits, e.g., patient falls
 - Challenge 3, forecast patient falls and evaluate the classification performance
 - □ Challenge 4, predict tremor dominance (TD) vs. posture instability and gait difficulty (PIGD).
- Results: model-free machine learning based techniques provide a more reliable clinical outcome forecasting, e.g., falls in Parkinson's patients, with classification accuracy of about 70-80%.

Gao, et al. SREP (2018)





Case-Studies – Parkinson's Disease

Method	асс	sens	spec	ppv	npv	lor	auc
Logistic Regression	0.728	0.537	0.855	0.710	0.736	1.920	0.774
Random Forests	<u>0.796</u>	<u>0.683</u>	<u>0.871</u>	<u>0.778</u>	<u>0.806</u>	<u>2.677</u>	<u>0.821</u>
AdaBoost	0.689	0.610	0.742	0.610	0.742	1.502	0.793
XGBoost	0.699	0.707	0.694	0.604	0.782	1.699	0.787
SVM	0.709	0.561	0.806	0.657	0.735	1.672	0.822
Neural Network	0.699	0.610	0.758	0.625	0.746	1.588	
Super Learner	0.738	0.683	0.774	0.667	0.787	1.999	

Results of binary fall/no-fall classification (5-fold CV) using top 10 selected features (gaitSpeed_Off, ABC, BMI, PIGD_score, X2.11, partII_sum, Attention, DGI, FOG_Q, H_and_Y_OFF)

Gao, et al. SREP (2018)



Open-Science & Collaborative Validation

End-to-end Big Data analytic protocol jointly processing complex imaging, genetics, clinical, demo data for assessing PD risk

- o Methods for rebalancing of imbalanced cohorts
- ML classification methods generating consistent and powerful phenotypic predictions
- Reproducible protocols for extraction of derived neuroimaging and genomics biomarkers for diagnostic forecasting

https://github.com/SOCR/PBDA





Case-Studies – ALS

- Identify predictive classifiers to detect, track and prognosticate the progression of ALS (in terms of clinical outcomes like ALSFRS and muscle function)
- Provide a decision tree prediction of adverse events based on subject phenotype and 0-3 month clinical assessment changes

Source	Sample Size/Data Type	Summary
ProAct Archive	Over 100 variables are recorded for all subjects including: <u>Demographics</u> : age, race, medical history, sex; <u>Clinical</u> data: <u>Amyotrophic Lateral Sclerosis</u> Functional Rating Scale (ALSFRS), adverse events, onset_delta, onset_site, drugs use (riluzole) The PRO-ACT training dataset contains clinical and lab test information of 8,635 patients. Information of 2,424 study subjects with valid gold standard ALSFRS slopes used for processing, modeling and analysis	The time points for all longitudinally varying data elements are aggregated into signature vectors. This facilitates the modeling and prediction of ALSFRS slope changes over the first three months (baseline to month 3)

Tang, et al. (2018), in review









Case-Studies – ALS – Dimensionality Reduction



Acknowledgments

Funding

NIH: P20 NR015331, U54 EB020406, P50 NS091856, P30 DK089503, P30AG053760, UL1TR002240 NSF: 1734853, 1636840, 1416953, 0716055, 1023115

The Elsie Andresen Fiske Research Fund

Collaborators

http://SOCR.umich.edu

- SOCR: Alexandr Kalinin, Selvam Palanimalai, Syed Husain, Matt Leventhal, Ashwini Khare, Rami Elkest, Abhishek Chowdhury, Patrick Tan, Gary Chan, Andy Foglia, Pratyush Pati, Brian Zhang, Juana Sanchez, Dennis Pearl, Kyle Siegrist, Rob Gould, Jingshu Xu, Nellie Ponarul, Ming Tang, Asiyah Lin, Nicolas Christou, Hanbo Sun, Tuo Wang. Simeone Marino
- LONIIN: Arthur Toga, Roger Woods, Jack Van Horn, Zhuowen Tu, Yonggang Shi, David Shattuck, Elizabeth Sowell, Katherine Narr, Anand Joshi, Shantanu Joshi, Paul Thompson, Luminita Vese, Stan Osher, Stefano Soatto, Seok Moon, Junning Li, Young Sung, Carl Kesselman, Fabio Macciardi, Federica Torri

 <u>UMich MIDAS/MNORC/AD/PD Centers</u>: Cathie Spino, Chuck Burant, Ben Hampstead, Stephen Goutman, Stephen Strobbe, Hiroko Dodge, Hank Paulson, Bill Dauer, Brian Athey



Acknowledgments

MIDAS Co-Directors

Brian Athey and Al Hero

MIDAS Education & Training Committee

Ivo Dinov HBBS/Bioinfo, Richard Gonzalez, ISR/PSY/LS&A, Eric Schwartz Ross & Kerby Shedden, Stats/LS&A

Inaugural Program Committee Members

H. V. Jagadish: Electrical Engineering and Computer Science, CoE Vijay Nair: Statistics & Industrial & Operations Engineering, LS&A/CoE George Alter: Institute for Social Research; History, LS&A Brian Athey: Computational Medicine and Bioinformatics, SoM Mike Cafarella: Computer Science and Engineering, CoE Ivo Dinov, Chair, Leadership and Effectiveness Science, Bioinformatics, SoN/SoM Karthik Duraisamy: Atmospheric, Oceanic, and Space Sciences August (Gus) Evrard: Physics; Astronomy, LS&A Anna Gilbert: Mathematics, LS&A Alfred Hero: Electrical Engineering and Computer Science; Biomedical Engineering, CoE Judy Jin: Industrial & Operations Engineering, CoE Carl Lagoze: School of Information Qiaozhu Mei: School of Information Christopher Miller: Astronomy, LS&A Dragomir Radev: School of Information; Computer Science and Engineering; Linguistics, CoE Stephen Smith: Ecology and Evolutionary Biology, LS&A Ambuj Tewari: Statistics; Computer Science and Engineering, LS&A Honglak Lee, Electrical Engineering and Computer Science, CoE Jeremy Taylor, Biostatistics, SPH



Michigan Institute for Data Science University of Michigan

www.MIDAS.umich.edu

Ivo Dinov dinov@umich.edu Open-ended discussion of educational challenges, research opportunities and infrastructure demands in data science

