Open Data Science & Predictive Health Analytics

Ivo D. Dinov

Statistics Online Computational Resource Health Behavior & Biological Sciences Computational Medicine & Bioinformatics Michigan Institute for Data Science

University of Michigan

http://SOCR.umich.edu

http://Predictive.Space

STRUCT OF NOTIONAL RESOURCE (SOCR)

Slides Online: "SOCR News"

Data Science and Predictive

Analytics

Outline

Driving biomedical & health challenges

□ Common characteristics of Big Neuroscience Data

Data science & predictive neuro-analytics methods

Case-studies

- □ Prenatal Exposure to Methamphetamine & Alcohol
- Parkinson's disease (PD)
- UK Biobank (Hands-on Demo)



Driving Biomedical/Health Challenges

□<u>Neurodegeneration</u>:

Structural Neuroimaging in Alzheimer's Disease illustrates the Big Data challenges in modeling complex neuroscientific data. 808 ADNI subjects, 3 groups: 200 subjects with Alzheimer's disease (AD), 383 subjects with mild cognitive impairment (MCI), and 225 asymptomatic normal controls (NC). The 80 neuroimaging biomarkers and 80 highly-associated SNPs.





Characteristics of Big Biomed Data

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

Tools
Harvesting and management of vast amounts of data
Wranglers for dealing with heterogeneous data
Tools for data harmonization and aggregation
Transfer and joint modeling of disparate elements
Macro to meso to micro scale observations
Techniques accounting for longitudinal patterns in the data
Reliable management of missing data

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, phenomics 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, et al. (2016) PMID:26918190



Data Science & Predictive Analytics

Data Science: an emerging extremely transdisciplinary field bridging between the theoretical, computational, experimental, and applied areas. Deals with enormous amounts of complex, incongruent and dynamic data from multiple sources. Aims to develop algorithms, methods, tools, and services capable of ingesting such datasets and supplying semi-automated decision support systems

Predictive Analytics: process utilizing advanced mathematical formulations, powerful statistical computing algorithms, efficient software tools, and distributed web-services to represent, interrogate, and interpret complex data. Aims to forecast trends, cluster patterns in the data, or prognosticate the process behavior either within the range or outside the range of the observed data (e.g., in the future, or at locations where data may not be available)



http://DSPA.predictive.space

Dinov, Springer (2018)



Data Science & Predictive Analytics

- Dimensionality Reduction
- Lazy Learning: Classification Using Nearest Neighbors
- Derived Probabilistic Learning: Classification Using Naive Bayes
- Decision Tree Divide and Conquer Classification
- Forecasting Numeric Data Using Regression Models
- Black Box Machine-Learning Methods: Neural Nets/Support Vector Machines
- □ Apriori Association Rules Learning
- L k-Means Clustering
- Model Performance Assessment
- Improving Model Performance
- Specialized Machine Learning Topics
- □ Variable/Feature Selection
- Regularized Linear Modeling and Controlled Variable Selection
- Big Longitudinal Data Analysis
- Natural Language Processing/Text Mining
- Prediction and Internal Statistical Cross Validation
- Deep Learning, Neural Networks

http://DSPA.predictive.space

Dinov, Springer (2018)





Case-Studies – Prenatal Exposure to Methamphetamine & Alcohol

- <u>Goals</u>: Examine the local brain effects of prenatal exposure to methamphetamine (MA)
 <u>Data</u>: structural magnetic resonance imaging (sMRI). Compared local brain volumes differed among 61 children (ages 5–15 years),
 - □ 21 with prenatal MA exposure,
 - □ 18 with concomitant prenatal alcohol exposure (the MAA group),
 - □ 13 with heavy prenatal alcohol but not MA exposure (ALC group), and
 - □ 27 unexposed controls.
- <u>Methods</u>: Brain morphometry (sMRI processing) & Discriminant analysis (prediction)
 Results:
 - Bilateral volume reductions in both exposure groups relative to controls in striatal and thalamic regions, right prefrontal and left occipitoparietal cortices.
 - MAA group had negative correlation between full-scale intelligence quotient (FSIQ) scores and caudate volume.
 - LDA prediction of group membership correctly classified 72% of participants.
- Conclusions: Striatal and limbic structures, known to be sites of neurotoxicity in adult MA abusers, may be more vulnerable to prenatal MA exposure than alcohol exposure; Severe striatal damage is associated with more severe cognitive deficit

Sowell, et al. JNeurosci (2010)



Case-Studies – Prenatal Exposure to Methamphetamine & Alcohol



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 controlled feature selection 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 validation + external out-of-bag validation
- Four specific challenges
 - 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)



Method	асс	sens	spec	рру	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	

Case-Studies – Parkinson's Disease

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 – General Populations

	20005	Ongoing characteri	stics Email access					
2	110007 Ongoing characteristics Newsletter communications, date sent					D LUC Disk and discriminate		
100	25780	Brain MRI Acc	quisition protocol phase.	Line and Line		UK BIODANK – disc	riminate	
100	12139	Brain MRI Bel	lieved safe to perform brain MRI scan			botwoon HC single	aand	
100	12188	Brain MRI Bra	ain MRI measurement completed			between no, singi	e anu	
100	12187	Brain MRI Bra	ain MRI measuring method			multiple comorbid	conditions	
100	12663	Brain MRI Rea	ason believed unsafe to perform brain I	MRI			oonaniono	
100	12704	Brain MRI Rea	ason brain MRI not completed			Predict likelihoods	of various	
100	12652	Brain MRI Rea	ason brain MRI not performed					
101	12292	Carotid ultrasound	Carotid ultrasound measurement	completed		developmental or a	aging	
101	12291	Carotid ultrasound	Carotid ultrasound measuring met	thod		diaardara	0 0	
101	20235	Carotid ultrasound	Carotid ultrasound results package	e		alsorders		
101	22672	Carotid ultrasound	Maximum carotid IMT (intima-me	edial thickness) a	t 1 20	Forecast cancer		
degree	25			_	_	i orecasi cancer		
101	22675	Carotid ultrasound	Maximum carotid IMT (intima-me	edial thickness) a	it 150			
degree	25)ata				
101	226/8	Carotid ultrasound	Maximum carotid IMT (intima-	Julu	Sam	nla Siza/Data Tuna	Summary	
				OUICO	Jan		Jullinuiv	
degree	25		5	ource	Jan	ipie Size/Data Type	Junnary	
degree 101	22681	Carotid ultrasound	S Maximum carotid IMT (intima-	ource	Dom	pographics: > 500K cases	The	
degree 101 degree	22681	Carotid ultrasound	S Maximum carotid IMT (intima-	ource	Dem	nographics: > 500K cases	The	
degree 101 degree 101	22681 22681 22671	Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med	ource	Dem	ographics: > 500K cases	The longitudinal	
degree 101 degree 101 101	22681 22681 22671 22674	Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med	JK	Derr Clini	nographics: > 500K cases cal data: > 4K features zing data: T1 resting-	The longitudinal archive of	
degree 101 degree 101 101 101	22681 22671 22674 22674 22677	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med	JK	Dem Clini Ima	nographics: > 500K cases ical data: > 4K features ging data: T1, resting-	The longitudinal archive of	
degree 101 degree 101 101 101 101	22681 22671 22674 22677 22680	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med	JK Biobank	Dem Clini Imag state	nographics: > 500K cases cal data: > 4K features ging data: T1, resting- e fMRI, task fMRI,	The longitudinal archive of the UK	
degree 101 degree 101 101 101 101 101	22681 22671 22674 22674 22677 22680 22670	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Minimum carotid IMT (intima-	JK Biobank	Dem Clini Imag state	nographics: > 500K cases cal data: > 4K features ging data: T1, resting- e fMRI, task fMRI, ELAIR, dMRI, SWI	The longitudinal archive of the UK population	
degree 101 degree 101 101 101 101 101 degree	22681 22671 22671 22674 22677 22680 22670 22670	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Minimum carotid IMT (intima- Minimum carotid IMT (intima-	JK Biobank	Dem Clini Imag state T2_F	nographics: > 500K cases ical data: > 4K features ging data: T1, resting- e fMRI, task fMRI, ELAIR, dMRI, SWI	The longitudinal archive of the UK population	
degree 101 degree 101 101 101 101 101 degree 101	22681 22671 22671 22674 22677 22680 22670 22670	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Minimum carotid IMT (intima- Minimum carotid IMT (intima-	JK Biobank	Dem Clini Imag state T2_F Gen	nographics: > 500K cases cal data: > 4K features ging data: T1, resting- e fMRI, task fMRI, FLAIR, dMRI, SWI etics data	The longitudinal archive of the UK population (NHS)	
degree 101 degree 101 101 101 101 101 degree 101	22681 22671 22674 22677 22680 22670 22670 22670 22673 22673	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Minimum carotid IMT (intima- Minimum carotid IMT (intima-	ource JK Biobank	Dem Clini Imag state T2_F Gen	nographics: > 500K cases cal data: > 4K features ging data: T1, resting- e fMRI, task fMRI, ELAIR, dMRI, SWI etics data	The longitudinal archive of the UK population (NHS)	
degree 101 degree 101 101 101 101 101 degree 101 degree	22681 22681 22671 22674 22677 22680 22670 22670 22673 22673 22673	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Minimum carotid IMT (intima- Minimum carotid IMT (intima- Minimum carotid IMT (intima-	ource JK Biobank dial thickness) at	Dem Clini Imag state T2_F Gen	http://www.ukbioban	The longitudinal archive of the UK population (NHS)	
degree 101 degree 101 101 101 101 101 degree 101 degree	22681 22681 22671 22674 22677 22680 22670 22670 22673 22676 22 22676	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Minimum carotid IMT (intima- Minimum carotid IMT (intima- Minimum carotid IMT (intima-	ource JK Biobank dial thickness) at	Dem Clini Imag state T2_F Gen t 210	http://www.ukbioban	The longitudinal archive of the UK population (NHS)	
degree 101 degree 101 101 101 101 101 101 degree 101 degree 101 degree	22681 22681 22671 22674 22677 22680 22670 22670 22673 22673 25 22676 25	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima- Minimum carotid IMT (intima- Minimum carotid IMT (intima- Minimum carotid IMT (intima- Minimum carotid IMT (intima-med	ource JK Biobank dial thickness) at dial thickness) at	Dem Clini Imag state T2_F Gen t 210 t 240	http://bd2k.org	The longitudinal archive of the UK population (NHS) tk.ac.uk	
degree 101 degree 101 101 101 101 101 degree 101 degree 101 degree 101	22681 22671 22674 22677 22680 22670 22670 22673 22673 22673 22675 22675 22679 22679 22679	Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima- Minimum carotid I	ource JK Biobank dial thickness) at dial thickness) at	Dem Clini Imag state T2_F Gen t 210 t 240	http://bd2k.org http://www.ukbioban	The longitudinal archive of the UK population (NHS) tk.ac.uk	
degree 101 degree 101 101 101 101 101 degree 101 degree 101 degree 101	22681 22671 22671 22677 22680 22670 22680 22670 22670 22673 22679 22679 22679 22679 22679	Carotid ultrasound Carotid ultrasound	S Maximum carotid IMT (intima- Mean carotid IMT (intima-med Mean carotid IMT (intima-med Mean carotid IMT (intima-med Minimum carotid IMT (intima- Minimum caroti	JK Biobank dial thickness) at dial thickness) at at 120 degrees t 150 degrees	Dem Clini Imag state T2_f Gen t 210 t 240	http://wyumi.ch/6wC	The longitudinal archive of the UK population (NHS) hk.ac.uk	

Case-Studies – UK Biobank

This Article Motivation Introduction	on Project Pipeline Unsupervised Learning TensorBoard Code Results Try-II-Now!	SOCR Home	
TensorBoard PROJECTOR		IRACTIVE 🔸 🔿 🕲	
DATA	Points: 9914 Dimension: 300 Selected 101 points	Show Al Isolete 101 Clear Data points selection	
1 tensor found Variable *		v Search ⊯ label =	5
Label y Color by label + label +		neighbors 0 • 10(8)	/it
label + No ignored label		distance COSINE EUCLIDEAN	L L
label + Tag selection as		Nearest points in the original space:	yo Yo
Load Download Label		1 0507 1 0504	-lt
Checkpoint: //ome/vegi/t_sne_webapp_vegi/ output_lifes/images.ckpt		a 0.676 1 0.675	♀ -
Metadata: /horw/vegi/Lane_velbapp_vegi/ output_Nes/dLlabels.tev		3 0.875 3 0.888	
TSNE PCA CUSTOM		2 0.699	b b
Dimension 20 (1) 30 Perplexity 0		1 0.712 1 0.713	at
Learning		2 0.718 1 0.723	a !
Supervise		2 0.725 2 0.730	
Iteration: 229		30.722 70.734	
How to use 1 SNE effectively.		BDOKMARKS (0) 🔍 🧄	
	http://myumi.ch/6wOgy		
http://c	nttp://mydninti/owogy	51	

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, Jerome Choi, Yuming Sun, Yi Zhao, Nina Zhou, Nicolas Christou, Hanbo Sun, Tuo Wang. Simeone Marino
- LONI/INI: 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
- UMich MIDAS/MNORC/AD/PD Centers: Cathie Spino, Chuck Burant, Ben Hampstead, Stephen Goutman, Stephen Strobbe, Hiroko Dodge, Hank Paulson, Bill Dauer, Brian Athey

