Predictive Big Data Analytics: Using Large, Complex, Heterogeneous, Incongruent, Multi-source & Incomplete Observations to Study Neurodegenerative Disorders

## lvo D. Dinov

Statistics Online Computational Resource Michigan Institute for Data Science School of Nursing University of Michigan

www.SOCR.umich.edu



## Outline

🗖 Big, Deep & Dark Data

Big Data Analytics

Applications to Neurodegenerative Disease

Data Dashboarding

Open Problem: Representation and Joint BD Analytics – Compressive Big Data Analytics (CBDA)



## Characteristics of Big Biomed 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



#### Mixture of quantitative & qualitative estimates

Dinov, et al., 2014

# Big, Deep & Dark Data Big: Size + Complexity + Variability +

Heterogeneous







## Data: Process & Model Representation

#### **Native Process**

#### **Model Representation**





Big Data	Information	Knowledge	Action
Raw Observations	Processed Data	Maps, Models	Actionable Decisions
Data Aggregation	Data Fusion	Causal Inference	Treatment Regimens
Data Scrubbing	Summary Stats	Networks, Analytics	Forecasts, Predictions
Semantic-Mapping	Derived Biomarkers	Linkages, Associations	Healthcare Outcomes



## Basic Core





## Big Data Analytics Resourceome



http://socr.umich.edu/docs/BD2K/BigDataResourceome.html



## Kryder's law: Exponential Growth of Data



Dinov, et al., 2014

## Big Data Challenges

Inference High-throughput Expeditive, Adaptive

#### Modeling

Constraints, Bio, Optimization, Computation

## Complexity

Volume, Heterogeneity

#### Representation

Incompleteness (space, time, measure)



#### Energy & Life-Span of Big Data

**Energy** encapsulates the holistic information content included in the data (exponentially increasing with time), which may often represent a significant portion of the joint distribution of the underlying healthcare process

Life-span refers to the relevance and importance of the Data (exponentially decreasing with time)



http://www.aaas.org/news/big-data-blog-part-v-interview-dr-ivo-dinov

## End-to-end Pipeline Workflow Solutions



A practical, hands-on guide to using the LONI Pipeline for neuroimaging analysis



Dinov, et al., 2014, Front. Neuroinform.;

Dinov, et al., Brain Imaging & Behavior, 2013



## End-to-end Pipeline Workflow Solutions



Dinov, et al., 2014, Front. Neuroinform.;

Dinov, et al., Brain Imaging & Behavior, 2013



Neurodegenerative Disease Applications: Imaging-Genetic Biomarker Interactions in <u>Alzheimer's Disease</u>

Investigate AD subjects (age 55 – 65) using Neuroimaging Initiative (ADNI) database to understand early-onset (EO) cognitive impairment using neuroimaging and genetics biomarkers

#### Data

EO-AD and EO-MCI (mild cognitive impairment) Derived 15 most impactful neuroimaging markers (out of 336 morphometry measures)

Obtained 20 most significant single nucleotide polymorphisms (SNPs) associated with specific neuroimaging biomarkers (out of 620K SNPs)

Approach

Global Shape Analysis (GSA) Pipeline workflow



Moon, Dinov et al., 2015, J Neuroimaging Moon, Dinov et al., 2015, Psych. Investig.

- Identified associations between neuroimaging phenotypes and genotypes for the EO cohort
- Overall most significant associations:
  - rs7718456 (Chr 15) and L\_hippocampus (volume)
  - rs7718456 and R\_hippocampus (volume)
- For the 27 EO-MCl's, most significant associations
  - rs6446443 (Chr 4, JAKMIP1 janus kinase and microtubule interacting protein 1 gene) and R\_superior\_frontal\_gyrus (volume)
  - rs17029131 (Chr 2) and L\_Precuneus (volume)



Moon, Dinov et al., 2015, Psych. Investig.









#### **QC** Protocol

A. QC Plink workflowB. Genetic association study

M

http://Pipeline.loni.usc.edu



Manhattan plot for all the single nucleotide polymorphisms (SNPs)

Neuroimaging phenotypes	p-value	Index	SNPs	Chromosome	p-value	Gene
L_cingulate_gyrus	0.0335	1	rs17029131	2	3.52E-06	
(Average mean curvature)						
L_gyrus_rectus (Surface area)	0.01728	2	rs1822144	2	2.28E-06	
R_cuneus (Surface area)	0.04203	3	rs6446443	4	6.68E-05	JAKMIP1 (janus kinase &
						microtubule interacting protein 1)
R_superior_frontal_gyrus (Volume)	0.03706	4	rs12506164	4	1.75E-05	
L_precentral_gyrus (Volume)	0.04125	5	rs7718456	5	3.36E-05	
L_precuneus (Volume)	0.0508	6	rs9377090	6	3.36E-05	
L_middle_occipital_gyrus	0.01805	7	rs2776932	10	2.20E-05	NRP1 (neuropilin1)
(Volume)						
R_superior_temporal_gyrus	0.03353	8	re4933672	10	648F-05	
(Volume)	0.05555	0	134755072	10	0.401-05	
L_hippocampus (Volume)	0.00067	9	rs11193270	10	3.52E-06	
R_hippocampus (Volume)	0.00539	10	rs11193272	10	3.52E-06	
R_precentral_gyrus (Shape index)	0.03411	11	rs11193274	10	3.52E-06	
R_precuneus (Shape index)	0.03186	12	rs12218153	10	3.52E-06	
L_cuneus (Shape index)	0.04952	13	rs1338956	10	2.20E-05	
R_inferior_occipital_gyrus	0.05037	14	rs1338025	10	2.20E-05	
(Curviness)						
R_putamen (Curviness)	0.03504	15	rs12101936	15	7.08E-06	
		16	rs16964473	19	3.53E-05	Intergenic
		17	rs12972537	19	6.14E-05	
		18	rs2212356	21	6.23E-05	
		19	rs2831165	21	6.68E-05	
		20	rs1266320	23	4.46E-06	

## **GWAS Imaging-Genetics Approach**

- **SNPs** 
  - E.g., C/T polymorphism
- Model





- **Phenotype:** Y<sub>i</sub> be the imaging-biomarker for *i*<sup>th</sup> subject
- Genotype: X<sub>i</sub> be the genotype *i*<sup>th</sup> subject at a particular:

$$SNP X_i = \begin{cases} 0, & BB \\ 1, & BA \\ 2, & AA \end{cases}$$

- **SOCR Multivariate Regression Models** 
  - $-Y_i = \beta_0 + \beta_1 X_i$
  - In general,  $Y_i = \sum_{k=0}^{K} \beta_k X_i^{(k)} + \varepsilon$
  - Stat analysis:  $\beta_k \neq 0$

Genotype-Phenotype Relation						
Pare	ents	O				
		В	А			
ę	В	BB	BA			
	A	BA	AA			



#### SNP-Neuroimaging interactions in Alzheimer's Disease



- Overall Associations
  - rs7718456 and L\_hippocampus (volume)
  - rs7718456 and R\_hippocampus (volume)
- EO-MCI associations
  - rs6446443
    - R\_superior\_frontal\_gyrus (volume)
- rs17029131 and L\_Precuneus (volume)



## Predictive Big Data Analytics in Parkinson's Disease

- A unique archive of Big Data: Parkinson's Progression Markers Initiative (PPMI). Data characteristics large size, incongruency, incompleteness, complexity, multiplicity of scales, and heterogeneity of information-generating sources
- **Approach** 
  - introduce methods for rebalancing imbalanced cohorts,
  - utilize a wide spectrum of classification methods to generate consistent and powerful phenotypic predictions,
  - generate reproducible machine-learning based classification that enables the reporting of model parameters and diagnostic forecasting based on new data.
- <u>Results</u> of machine-learning based classification show significant power to predict Parkinson's disease in the PPMI subjects (consistent accuracy, sensitivity, and specificity exceeding 96%, confirmed using statistical n-fold cross-validation).
- Model-free machine learning-based classification methods (e.g., adaptive boosting, support vector machines) outperform model-based techniques in terms of predictive precision and reliability (e.g., forecasting patient diagnosis).



#### Predictive Big Data Analytics using Large, Complex, Incongruent, Heterogeneous, Multi-source & Incomplete Observations

#### A Big Data Study of Parkinson's Disease





## SOCR Big Data Dashboard http://socr.umich.edu/HTML5/Dashboard

- Web-service combining and integrating multi-source socioeconomic and medical datasets
- Big data analytic processing
- Interface for exploratory navigation, manipulation and visualization
- Adding/removing of visual queries and interactive exploration of multivariate associations

Powerful HTML5 technology enabling mobile on-demand computing

Husain, et al., 2015, J Big Data



#### SOCR Dashboard (Exploratory Big Data Analytics): Data Fusion



#### SOCR Dashboard (Exploratory Big Data Analytics): Data QC





http://socr.umich.edu/HTML5/Dashboard

#### SOCR Dashboard (Exploratory Big Data Analytics): Associations



Open Problem: Representation and Joint BD Analytics

One Approach: Compressive Big Data Analytics (CBDA)

![](_page_27_Picture_2.jpeg)

#### **Motivation!** Big Data Analytics – Compressive Sensing

- There is currently no established analytical foundation for systematic representation of Big Data that facilitates the handling of data complexities and at the same time enables joint modeling, information extraction, high-throughput and adaptive scientific inference
- One idea is Compressive Big Data Analytics (CBDA), which borrows some of the compelling ideas for representation, reconstruction, recovery and data denoising recently developed for compressive sensing
- In compressive sensing, a sparse (incomplete) data is observed and one looks for a high-fidelity estimation of the complete dataset.
   Sparse data (or signals) can be described as observations with a small support, i.e., small magnitude according to the zero-norm

![](_page_28_Picture_4.jpeg)

#### Big Data Analytics – Compressive Sensing

![](_page_29_Figure_1.jpeg)

- The foundation for Compressive Big Data Analytics (CBDA) involves
  - Iteratively generating random (sub)samples from the Big Data collection.
  - Then, using classical techniques to obtain model-based or non-parametric inference based on the sample.
  - Next, compute likelihood estimates (e.g., probability values quantifying effects, relations, sizes)
  - Repeat the process continues iteratively.
- Repeating the (re)sampling and inference steps many times (with or without using the results of previous iterations as priors for subsequent steps).

![](_page_30_Picture_7.jpeg)

- Finally, bootstrapping techniques may be employed to quantify joint probabilities, estimate likelihoods, predict associations, identify trends, forecast future outcomes, or assess accuracy of findings.
- The goals of compressive sensing and compressive big data analytics are different.
  - CS aims to obtain a stochastic estimate of a complete dataset using sparsely sampled incomplete observations.
  - CBDA attempts to obtain a quantitative joint inference characterizing likelihoods, tendencies, prognoses, or relationships.
  - However, a common objective of both problem formulations is the optimality (e.g., reliability, consistency) of their corresponding estimates.

- Suppose we represent (observed) Big Data as a large matrix  $Y \in R^{n \times t}$ , where n= sample size (instances) and t = elements (e.g., time, space, measurements, etc.)
- To formulate the problem in an <u>analytical framework</u>, let's assume  $L \in \mathbb{R}^{n \times t}$  is a low rank matrix representing the mean or background data features,  $D \in \mathbb{R}^{n \times m}$  is a (known or unknown) design or dictionary matrix,  $S \in \mathbb{R}^{m \times t}$  is a sparce parameter matrix with small support  $(supp(S) \ll m \times t), E \in \mathbb{R}^{n \times t}$  denote the model error term, and  $\Lambda_{\Omega}(.)$  be a sampling operator generating incomplete data over the indexing pairs of instances and data elements

 $\Omega \subseteq \{1, 2, \dots, n\} \times \{1, 2, \dots, t\}$ 

• In this generalized model setting, the problem formulation involves estimation of *L*, *S* (and *D*, if it is unknown), according to this model representation:  $\Lambda_{\Omega}(Y) = \Lambda_{\Omega}(L + DS + E)$  (2)

![](_page_32_Picture_5.jpeg)

- Having quick, reliable and efficient estimates of L, S and D would allow us to make inference, compute likelihoods (e.g., p-values), predict trends, forecast outcomes, and adapt the model to obtain revised inference using new data
- When D is known, the model in equation (2) is jointly convex for L and S, and there exist iterative solvers based on sub-gradient recursion (e.g., alternating direction method of multipliers)
- However, in practice, the size of Big Datasets presents significant computational problems, related to slow algorithm convergence, for estimating these components that are critical for the final study inference

![](_page_33_Picture_4.jpeg)

- One strategy for tackling this optimization problem is to use a random Gaussian sub-sampling matrix  $A_{m \times n}$  (much like in the compressive sensing protocol) to reduce the rank of the observed data  $(Y_{m \times l}, where (m, l) \in \Omega)$  and then solve the minimization using least squares
- This *partitioning* of the difficult general problem into smaller chunks has several advantages. It reduces the hardware and computational burden, enables algorithmic parallelization of the global solution, and ensures feasibility of the analytical results
- Because of the stochastic nature of the index sampling, this approach may have desirable analytical properties like predictable asymptotic behavior, limited error bounds, estimates' optimality and consistency characteristics

![](_page_34_Picture_4.jpeg)

![](_page_35_Figure_1.jpeg)

Data Structure (Representation)

![](_page_35_Picture_3.jpeg)

#### Sample Data (Instance)

![](_page_35_Picture_5.jpeg)

• One can design an algorithm that searches and keeps only the most informative data elements by requiring that the derived estimates represent optimal approximations to y within a specific sampling index subspace  $\{(m, l)\} \subseteq \Omega$ 

 We want to investigate if CBDA inference estimates can be shown to obey error bounds similar to the upper bound results of point imbedding's in high-dimensions (e.g., Johnson-Lindenstrauss lemma) or the restricted isometry property

![](_page_36_Picture_3.jpeg)

• The Johnson-Lindenstrauss lemma guarantees that for any  $0 < \varepsilon < 1$ , a set of points  $\{P_k\}_1^K \in R^n$  can be linearly embedded  $(\Psi: R^n \to R^{n'})$  into  $\{\Psi(P_k) = P_k'\}_1^K \in R^{n'}$ , for  $\forall n' \ge 4\left(\frac{ln(K)}{\frac{\epsilon^2}{2}-\frac{\epsilon^3}{2}}\right)$ ,

almost preserving their pairwise distances, i.e.,

$$(1-\epsilon) \|P_i - P_j\|_2^2 \le \|P_i' - P_j'\|_2^2 \le (1+\epsilon) \|P_i - P_j\|_2^2$$

The restricted isometry property ensures that if δ<sub>2k</sub> < √2 − 1 and the estimate x̂ = arg min<sub>z:Az=y</sub> ||z||<sub>1</sub>, where A<sub>m×n</sub> satisfies property (1), then the data reconstruction is reasonable, i.e.,

$$\|\hat{x} - x\|_2 \le C_0 \frac{\sigma_k(x)_1}{\sqrt{k}}$$

 Can we develop iterative space-partitioning CBDA algorithms that either converge to a fix point or generate estimates that are *close* to their corresponding inferential parameters?

#### Acknowledgments

#### Funding

#### NIH: P50 NS091856, P30 DK089503, P20 NR015331, U54 EB020406 NSF: 1416953, 0716055 1023115 http://SOCR.umich.edu

#### **Collaborators**

- 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
- 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, Fabio Macciardi, Federica Torri, Carl Kesselman,
- UMich AD/PD Centers: Cathie Spino, Ben Hampstead, Bill Dauer

![](_page_38_Picture_7.jpeg)