

# Big Brain Data & Predictive Analytics

Ivo D. Dinov

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

University of Michigan

[www.SOCR.umich.edu](http://www.SOCR.umich.edu)



SCHOOL OF NURSING

STATISTICS ONLINE COMPUTATIONAL RESOURCE (SOCR)

UNIVERSITY OF MICHIGAN

## Outline

- ❑ AMIA Biomedical Imaging Working Group Driving Challenge:  
*Does Integrative Data Analytics on Biomedical Imaging  
Bring Us Closer to Precision Medicine?*
- ❑ Common characteristics of Big Brain Data
- ❑ Data science & predictive analytics
- ❑ Case-studies
- ❑ Applications to Neurodegenerative Disease
- ❑ Data Dashboarding
- ❑ Compressive Big Data Analytics (CBDA)



## Integrative Data Analytics ↔ Precision Medicine

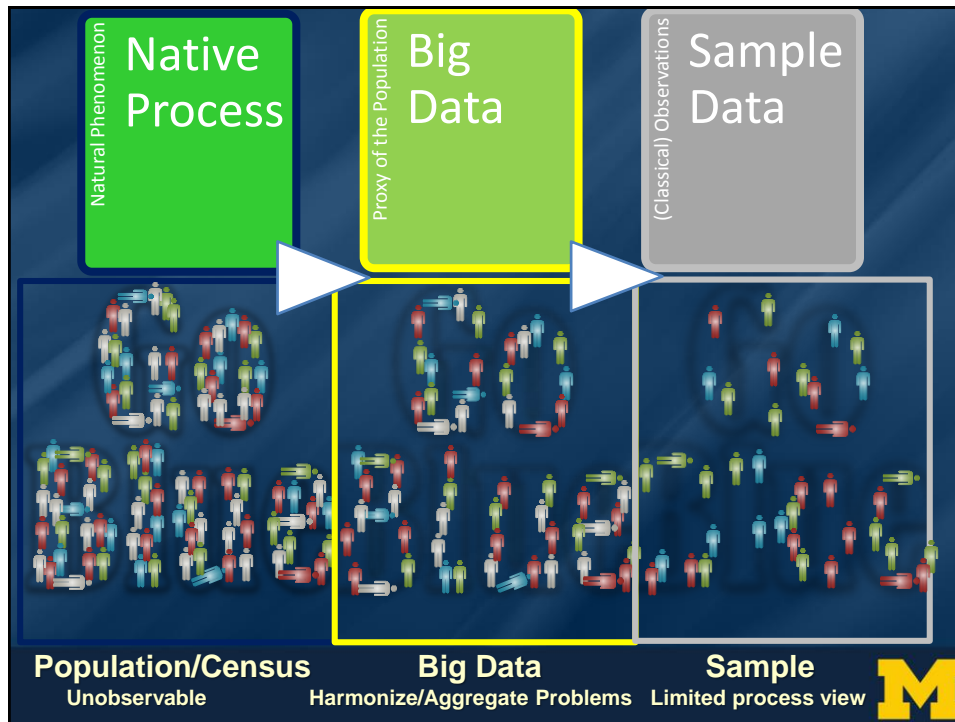
### □ **Neurodegeneration:**

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.



## Integrative Data Analytics ↔ Precision Medicine

- **Information Complexity** – large, incongruent, time-varying data
- **Precision Medicine** – customized medical decisions, clinical practice, treatments, or healthcare products to individual patients
- **Individual vs. Population Studies** – inductive (discriminative) vs. deductive (generative) models for clinical decision support
- **Tools** – molecular diagnostics, imaging, clinical, wearables, analytics, ...



## Characteristics of Big Biomed Data

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

### Big Bio Data Dimensions

### Tools

|              |   |
|--------------|---|
| Size         | Harvesting and management of vast amounts of data |
| Complexity   | Wranglers for dealing with heterogeneous data     |
| Incongruency | Tools for data harmonization and aggregation      |
| Multi-source | Transfer and joint modeling of disparate elements |
| Multi-scale  | Macro to meso to micro scale observations         |
| Incomplete   | 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

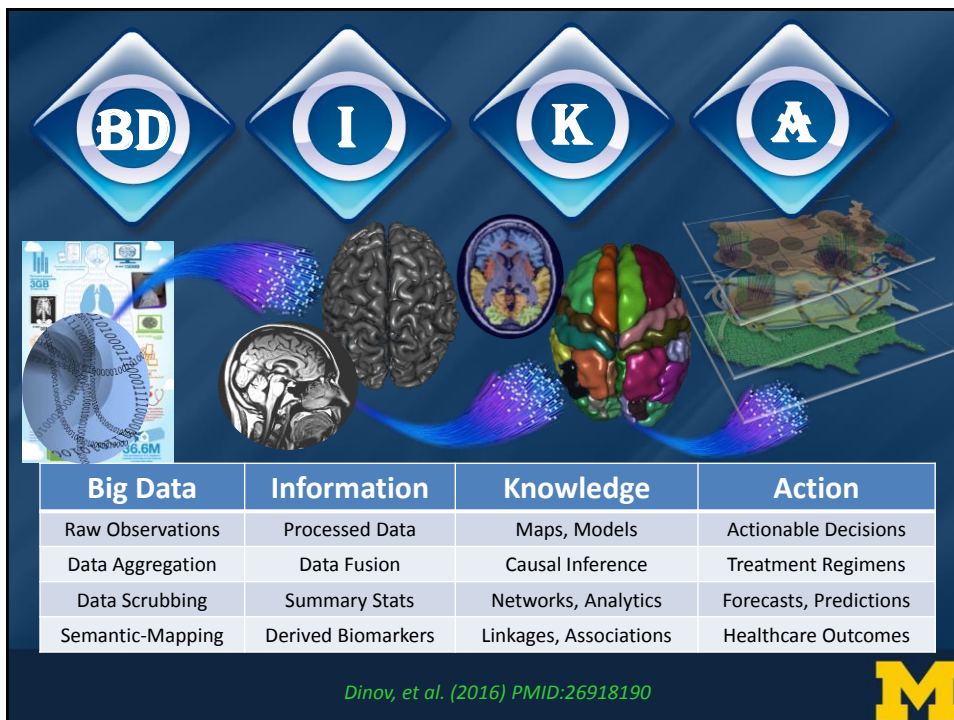
Dinov, *et al.* (2016) PMID:26918190



# Data science & predictive analytics

- ❑ **Data science:** an emerging extremely transdisciplinary field - bridging between the theoretical, computational, experimental, and biosocial 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:** utilizing advanced mathematical formulations, powerful statistical computing algorithms, efficient software tools and 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>



## Case-Studies – ALS

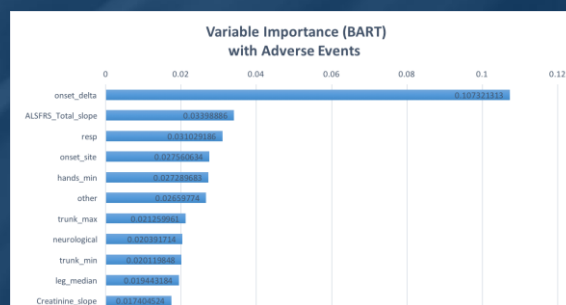
- Identify predictive classifiers to detect, track and prognosticate the progression of ALS (in terms of clinical outcomes like ALSFRS and muscle function)
- Provide a decision tree prediction of adverse events based on subject phenotype and 0-3 month clinical assessment changes

| Data 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 data</u> : Amyotrophic Lateral Sclerosis Functional Rating Scale (ALSFRS), adverse events, onset_delta, onset_site, drugs use (riluzole)<br>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) |



## Case-Studies – ALS

- Detect, track and prognosticate the progression of ALS
- Predict of adverse events based on subject phenotype and 0-3 month clinical assessment changes



| Methods     | Linear Regression | Random Forest | BART         | SuperLearner |
|-------------|-------------------|---------------|--------------|--------------|
| R-squared   | 0.081             | 0.174         | <b>0.225</b> | 0.178        |
| RMSE        | 0.619             | 0.587         | <b>0.568</b> | 0.585        |
| Correlation | 0.298             | 0.434         | <b>0.485</b> | 0.447        |





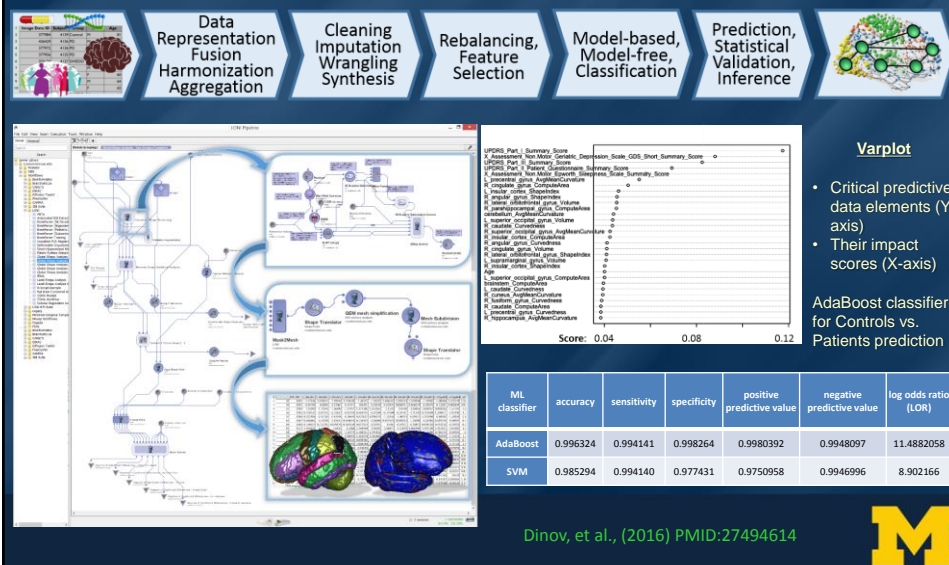
# Case-Studies – Parkinson's Disease

- Predict the clinical diagnosis of patients using all available data (with and without the UPDRS clinical assessment, which is the basis of the clinical diagnosis by a physician)
- Compute derived neuroimaging and genetics biomarkers that can be used to model the disease progression and provide automated clinical decisions support
- Generate decision trees for numeric and categorical responses (representing clinically relevant outcome variables) that can be used to suggest an appropriate course of treatment for specific clinical phenotypes

| Data Source  | Sample Size/Data Type  | Summary   |
|--------------|--|---|
| PPMI Archive | <p>Demographics: age, medical history, sex.</p> <p><u>Clinical</u> data: physical, verbal learning and language, neurological and olfactory, UPSIT, UPDRS scores, ADL, GDS-15, ...</p> <p><u>Imaging</u> data: structural MRI.</p> <p><u>Genetics</u> data: APOE genotypes e2/e3</p> <p><u>Cohorts</u>: Group 1 = {PD Subjects}, <math>N_1 = 263</math>; Group 2 = {PD Subjects with Scans without Evidence of a Dopaminergic Deficit (SWEDD)}, <math>N_2 = 40</math>; Group 3 = {Control Subjects}, <math>N_3 = 127</math>.</p> | <p>The longitudinal PPMI dataset including clinical, biological and imaging data (screening, baseline, 12, 24, and 48 month follow-ups) may be used conduct model-based predictions as well as model-free classification and forecasting analyses</p> |



## Predictive Big Data Analytics: Applications to Parkinson's Disease



## Case-Studies – General Populations

|     |        |                         |  |
|-----|--------|-------------------------|--|
| 2   | 20005  | Ongoing characteristics | Email access   |
| 2   | 110007 | Ongoing characteristics | Newsletter communications, date sent                         |
| 100 | 25780  | Brain MRI               | Acquisition protocol phase.                                  |
| 100 | 12139  | Brain MRI               | Believed safe to perform brain MRI scan                      |
| 100 | 12188  | Brain MRI               | Brain MRI measurement completed                              |
| 100 | 12187  | Brain MRI               | Brain MRI measuring method                                   |
| 100 | 12663  | Brain MRI               | Reason believed unsafe to perform brain MRI                  |
| 100 | 12704  | Brain MRI               | Reason brain MRI not completed                               |
| 100 | 12652  | Brain MRI               | Reason brain MRI not performed                               |
| 101 | 12292  | Carotid ultrasound      | Carotid ultrasound measurement completed                     |
| 101 | 12291  | Carotid ultrasound      | Carotid ultrasound measuring method                          |
| 101 | 20235  | Carotid ultrasound      | Carotid ultrasound results package                           |
| 101 | 22672  | Carotid ultrasound      | Maximum carotid IMT (intima-medial thickness) at 120 degrees |
| 101 | 22675  | Carotid ultrasound      | Maximum carotid IMT (intima-medial thickness) at 150 degrees |
| 101 | 22678  | Carotid ultrasound      | Maximum carotid IMT (intima-medial thickness) at 210 degrees |
| 101 | 22681  | Carotid ultrasound      | Maximum carotid IMT (intima-medial thickness) at 240 degrees |
| 101 | 22671  | Carotid ultrasound      | Mean carotid IMT (intima-medial thickness) at 120 degrees    |
| 101 | 22674  | Carotid ultrasound      | Mean carotid IMT (intima-medial thickness) at 150 degrees    |
| 101 | 22677  | Carotid ultrasound      | Mean carotid IMT (intima-medial thickness) at 210 degrees    |
| 101 | 22680  | Carotid ultrasound      | Mean carotid IMT (intima-medial thickness) at 240 degrees    |
| 101 | 22670  | Carotid ultrasound      | Minimum carotid IMT (intima-medial thickness) at 120 degrees |
| 101 | 22673  | Carotid ultrasound      | Minimum carotid IMT (intima-medial thickness) at 150 degrees |
| 101 | 22676  | Carotid ultrasound      | Minimum carotid IMT (intima-medial thickness) at 210 degrees |
| 101 | 22679  | Carotid ultrasound      | Minimum carotid IMT (intima-medial thickness) at 240 degrees |
| 101 | 22682  | Carotid ultrasound      | Quality control indicator for IMT at 120 degrees             |
| 101 | 22683  | Carotid ultrasound      | Quality control indicator for IMT at 150 degrees             |
| 101 | 22684  | Carotid ultrasound      | Quality control indicator for IMT at 210 degrees             |
| 101 | 22685  | Carotid ultrasound      | Quality control indicator for IMT at 240 degrees             |

- UK Biobank – discriminate between HC, single and multiple comorbid conditions
- Predict likelihoods of various developmental or aging disorders
- Forecast cancer

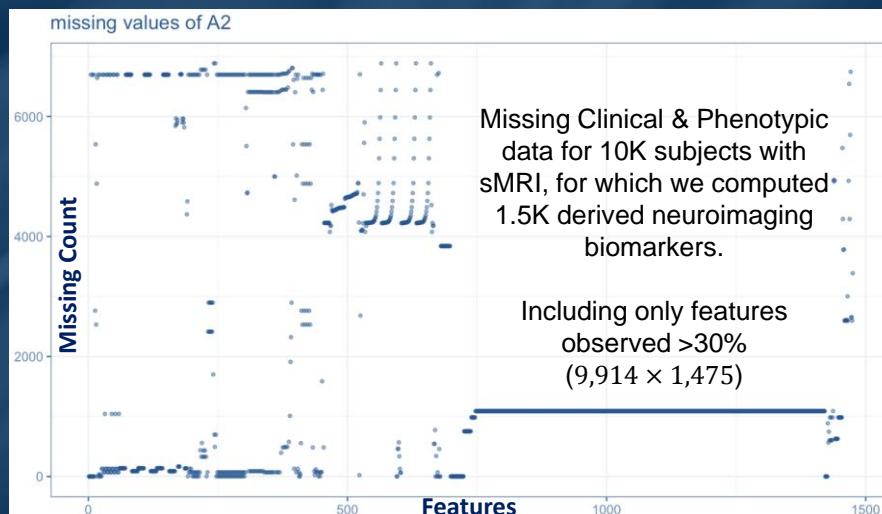
| Data Source | Sample Size/Data Type   | Summary   |
|-------------|---|---|
| UK Biobank  | <b>Demographics:</b> > 500K cases<br><b>Clinical data:</b> > 4K features<br><b>Imaging data:</b> T1, resting-state fMRI, task fMRI, T2_FLAIR, dMRI, SWI<br><b>Genetics data</b> | The longitudinal archive of the UK population (NHS) |

<http://www.ukbiobank.ac.uk>

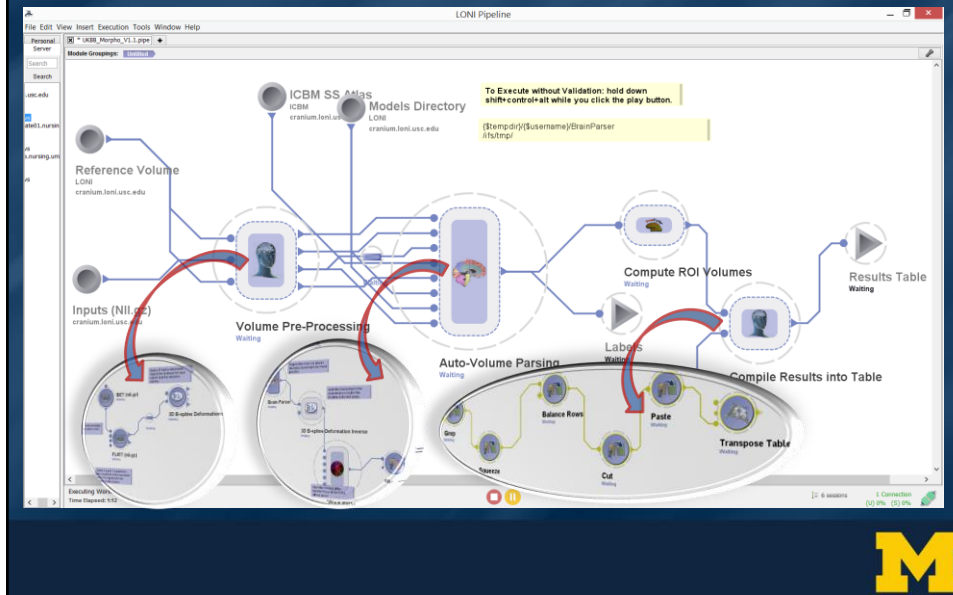
<http://bd2k.org>



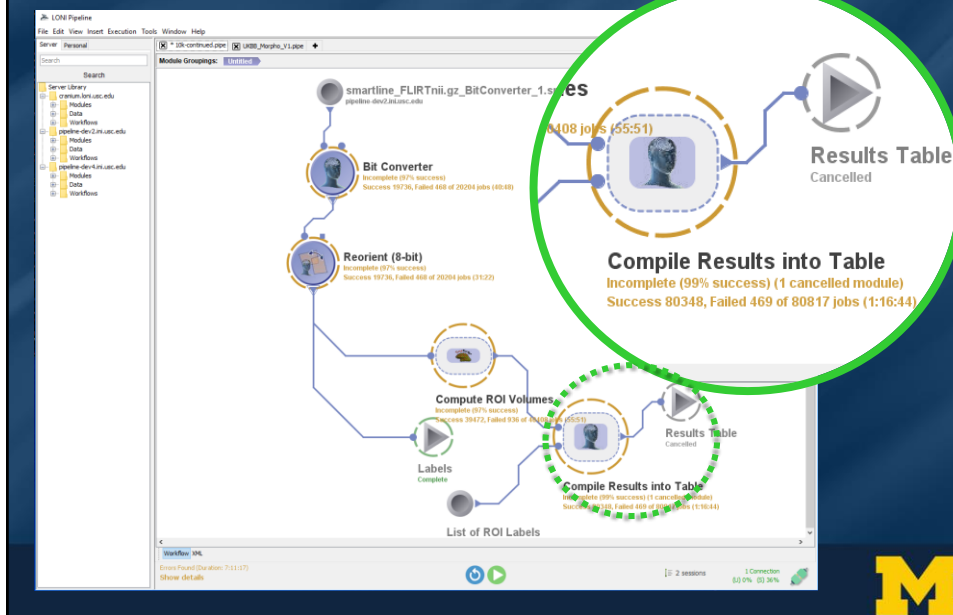
## Case-Studies – UK Biobank - Complexities



## Case-Studies – UK Biobank – NI Biomarkers



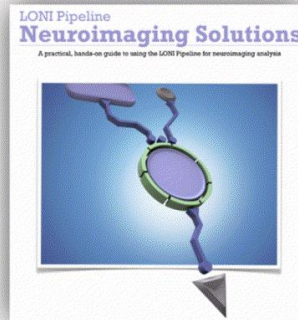
## Case-Studies – UK Biobank – Successes/Failures





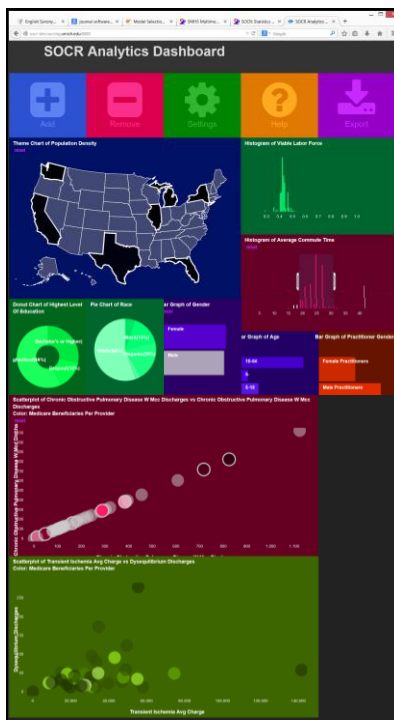


# End-to-end Pipeline Workflow Solutions



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

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



## SOCR Big Data Dashboard

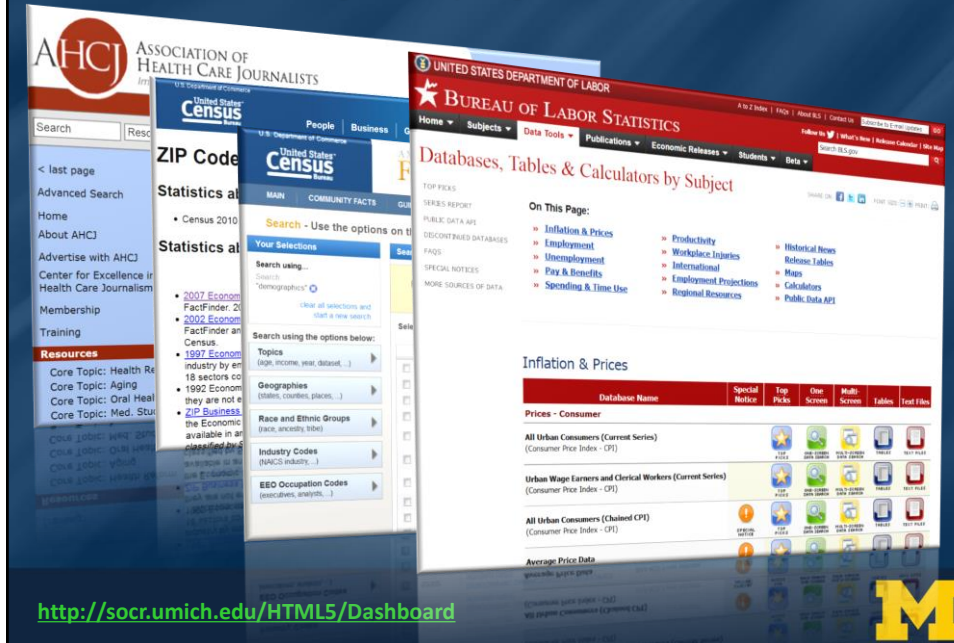
<http://socr.umich.edu/HTML5/Dashboard>

- ❑ 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

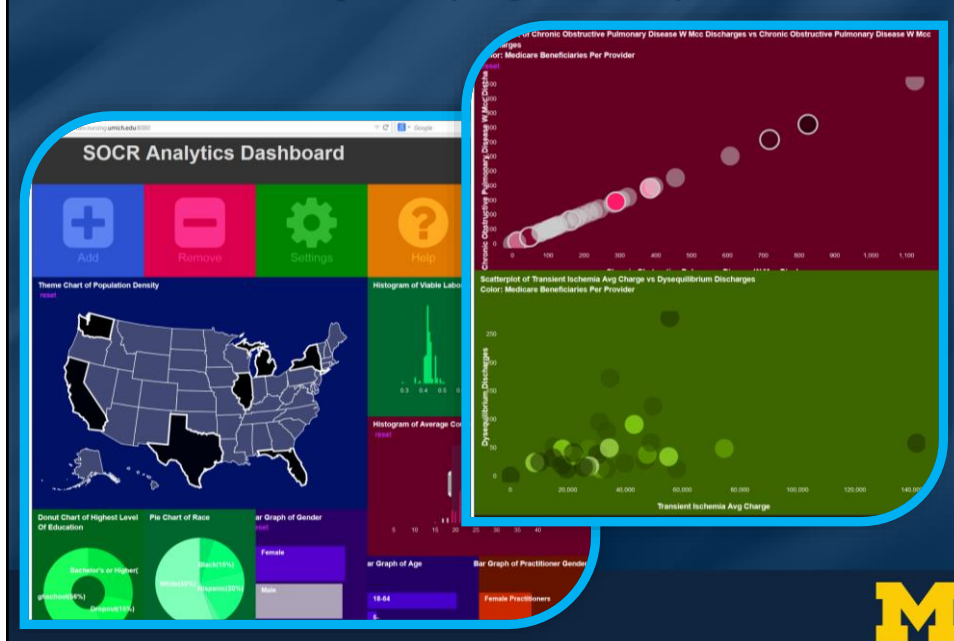
Husain, *et al.*, 2015, *PMID:26236573*



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



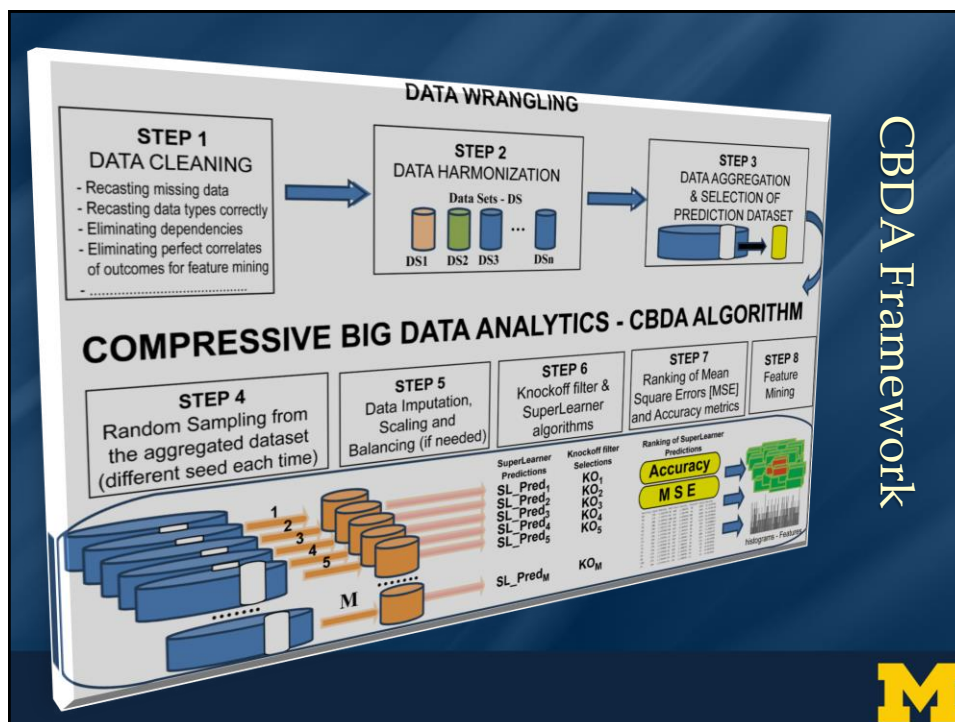
## SOCR Dashboard (Exploratory Big Data Analytics): Associations



# Compressive Big Data Analytics (CBDA)

- Foundations for Compressive Big Data Analytics (CBDA)
  - Iteratively generate 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 until a criterion is met – the (re)sampling and inference steps many times (with or without using the results of previous iterations as priors for subsequent steps)

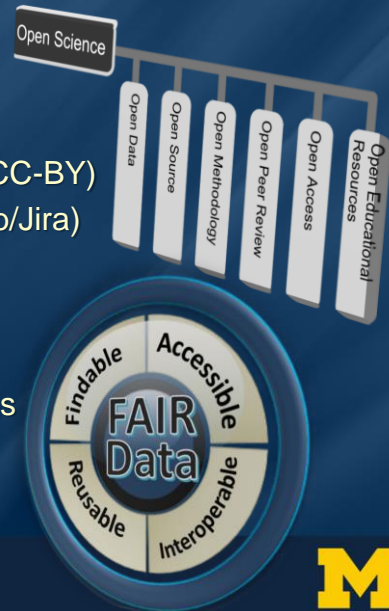
Dinov, 2016, PMID: 26998309





## FAIR Data & Open-Science Principles

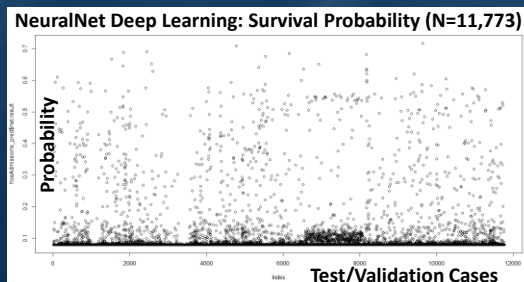
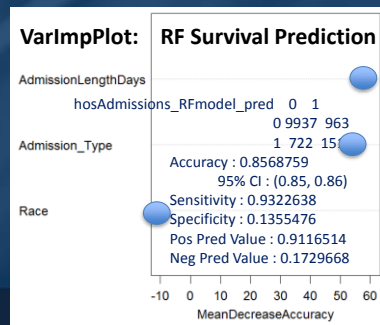
- ☐ Share resources
- ☐ Collaborate
- ☐ Permissive licenses (e.g., LGPL/CC-BY)
- ☐ Project management (e.g., GitHub/Jira)
- ☐ Open-access pubs
- ☐ Public-private partnerships
- ☐ Co-mentoring of trainees
- ☐ Effective transdisciplinary methods
- ☐ Resource Interoperability
- ☐ Result Reproducibility



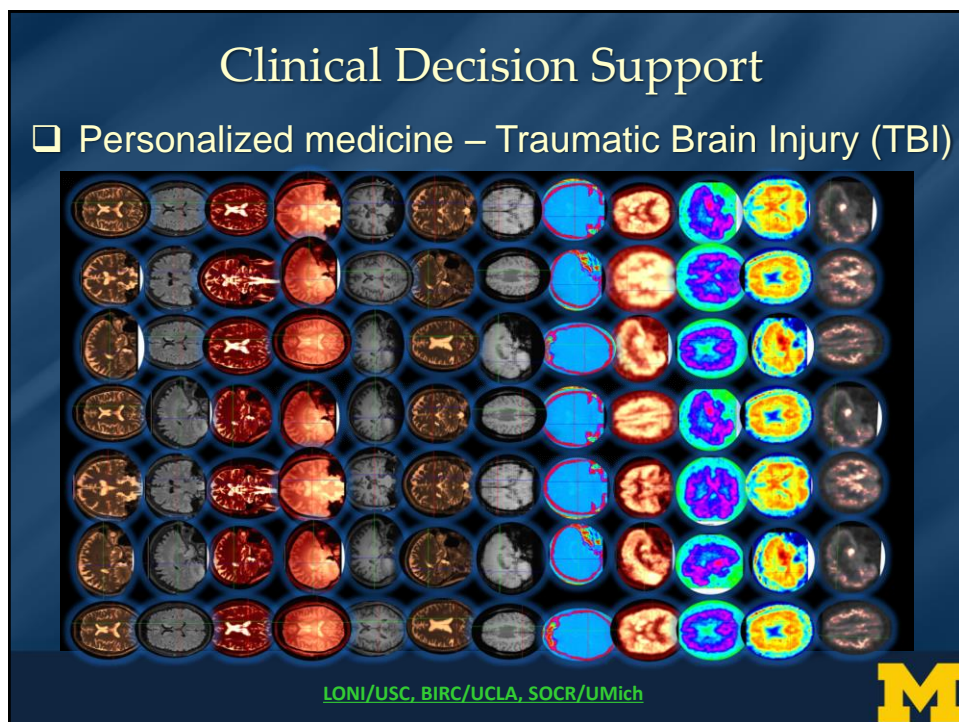
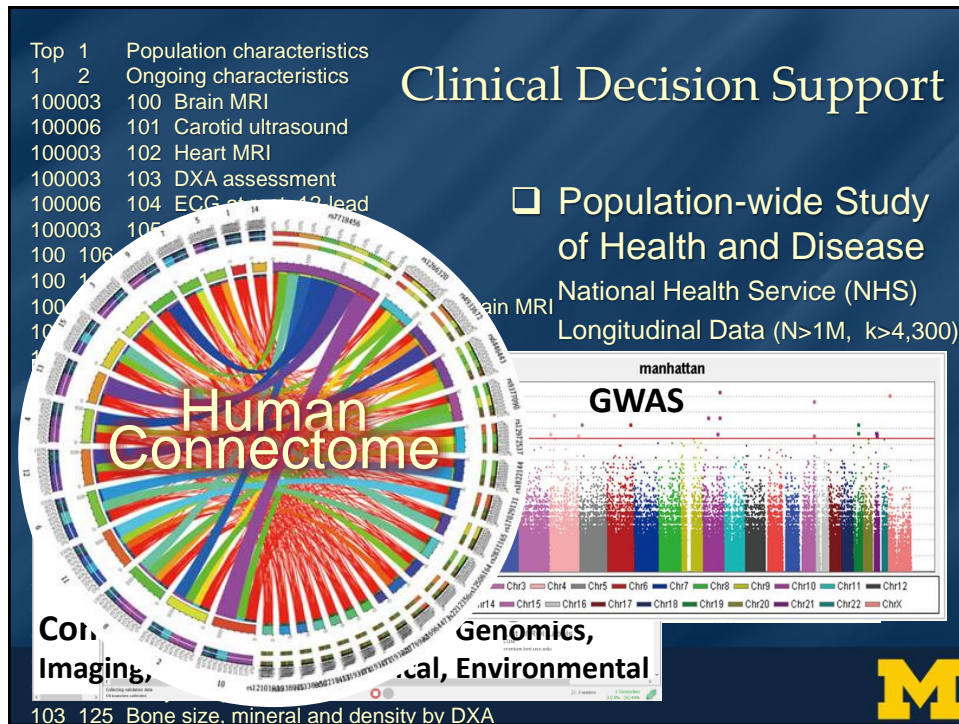
## Clinical Decision Support

### ☐ Hospital Admissions

- ☐ Survival Inference and Clinical outcome forecasting using hospital admissions data (N~60K and k=9):
  - ☐ Admission\_Length: Duration of hospital stay (in days)
  - ☐ Death: Indicator of Death (1) or survival (0)
  - ☐ ...
  - ☐ Demographics & (Human-labeled) Diagnoses







## Acknowledgments

Slides available online:

Google "SOCR News"

### Funding

NIH: P20 NR015331, U54 EB020406, P50 NS091856, P30 DK089503, P30AG053760

NSF: 1734853, 1636840, 1416953, 0716055, 1023115

The Elsie Andresen Fiske Research Fund

<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, Hanbo Sun, Tuo Wang
- **LONI/NI:** 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

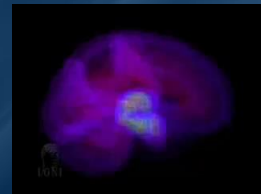


SOCR  
STATISTICS



## Demo(s)?

- Complex DB Search, retrieval (IDA)
- Multidimensional data visualization (MotionCharts, BrainViewer, R)
- Distributed high-throughput pipeline workflow computing
- SOCRAT Framework
- Data Dashboard
- Education and Training Resources
  - Probability and Statistics Ebook (EBook)
  - Scientific Methods for Health Sciences (SMHS)
  - Data Science and Predictive Analytics (DSPA) MOOC
  - SOCR Tools (distribution calculators, charts, modeler, analyses, experiments)
- Compressive Big Data Analytics (CBDA)



Demo

