# Big Brain Data Science and 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

www.SOCR.umich.edu

STATISTICS ONLINE COMPUTATIONAL RESOURCE (SOCR)

# Outline

- Driving biomedical & health challenges
- Common characteristics of Big Brain Data
- Data science & predictive analytics
- Case-studies
  - Applications to Neurodegenerative Disease
  - Data Dashboarding
- □ Compressive Big Data Analytics (CBDA)

# 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.



| Driving biomedical/health challenges   |
|--|
| <ul> <li>Phenotype-Genotype-Environmental</li> <li>Phenotype:</li> <li>Outward manifestations of the form, shape, or characteristics of specific individuals, cohorts, or morbid conditions</li> <li>Nature/Genome:</li> <li>Read the genome →</li> <li>protein synthesis →</li> <li>2 basic functions →</li> <li>structural proteins determine <i>physical traits</i> or <i>functional traits</i> of protein enzymes catalyzing chemical reactions</li> <li>Nurture/Environment:</li> <li>Ambient environment directly effects Pheno + Geno →</li> <li>Provides raw materials needed for the synthetic processes controlled by Nature/Genes →</li> <li>Organisms either synthesize of obtain amino acids with their diet</li> </ul> |
|  |



# Characteristics of Big Biomed Data

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

| Big Bio Data<br>Dimensions | Tools  |
|----------------------------|--|
| Size                       | Harvesting and management of vast amounts of data    |
| Complexity                 | Wranglers for dealing with<br>heterogeneous data     |
| Incongruency               | Tools for data harmonization and aggregation         |
| Multi-source               | Transfer and joint modeling of<br>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: process 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

Dinov, Springer (2018)





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





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



## **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



## o Goals

- 1) Harmonize and aggregate complex, multisource, and multi-site PD data
- Identify highly predictive features associated with specific clinical traits, e.g., falls
- Forecast falls using ML techniques and validate using statistical methods





## PD Patient Falls – Predictive BD Analytics

## **Results**

Feature selection for the Udall data using RF and KO.

7 common features selected by both methods: MDS\_PIGD, gaitSpeed\_Off, MOT\_EDL, NON\_MOTOR\_EDL, walk, pos\_stab

### Binary classification of falls/no-fall

(5-fold CV) using top 10 selected features (gaitSpeed\_Off, ABC, BMI, PIGD\_score, X2.11, partll\_sum, Attention, DGI, FOG\_Q, H\_and\_Y\_OFF)

Method

AdaBoost

XGBoost

Neural Network

Super Learner

Logit

|                           |   | Rar          | ndom Forest | s (RF)              |        |               |           | Knockoff (K | 0)    |
|---------------------------|---|--------------|-------------|---------------------|--------|---------------|-----------|-------------|-------|
| Features                  |   |              |             | Free                | quency |               | Features  | Frequency   |       |
|                           |   | MDS_         | PIGD        |                     | 0      | .888          |           | nx_smoke    | 0.764 |
|                           | gaitSpeed_Off   |              |             |                     |        | 0.860 high_bp |           |             | 0.751 |
| R_middle_temporal_gyrus   |   |              |             | 0                   | .662   |               | 0.718     |             |       |
| R_inferior_temporal_gyrus |   |              |             | 0                   | .618   | Ν             | /IDS_PIGD | 0.672       |       |
|                           |   | Cauda        | te_DA       |                     | 0      | .554          | SL        | EEP_APNEA   | 0.602 |
|                           |   | Striatu      | m_DA        |                     | 0      | .534          |           | head_inj    | 0.598 |
|                           |   | MOT          | EDL         |                     | 0      | .516          | S         | LEEP_RBD    | 0.552 |
|                           | time_upgo<br>L_middle_temporal_gyrus<br>NON_MOTOR_EDL |              | 0           | .494                |        | out_bed       | 0.515     |             |       |
|                           |   |              | 0           | 0.436 gaitSpeed_Off |        | 0.502         |           |             |       |
|                           |   |              |             | 0                   | .418   |               | НҮ        | 0.477       |       |
|                           |   | UPS          | IT40        |                     | 0      | .410          | NON       | _MOTOR_EDL  | 0.440 |
|                           |   | Putam        | en_DA       |                     | 0      | .408          |           | hal_psy     | 0.415 |
|                           | R_  | middle_orbit | ofrontal_gy | rus                 | 0      | .364          |           | chair       | 0.415 |
|                           | spec  | рру          | npv         |                     | lor    | 354           |           | pos_stab    | 0.407 |
|                           | 0.855   | 0.710        | 0.736       |                     | 1.920  | 336           | Ci        | audate_DA   | 0.403 |
| 0                         | .871  | 0.778        | 0.806       | 2                   | 2.677  | 324           |           | MOT_EDL     | 0.398 |
|                           | 0.742   | 0.610        | 0.742       |                     | 1.502  | 322           |           | gait        | 0.374 |
|                           | 0.694   | 0.604        | 0.782       |                     | 1.699  | 320           |           | gender      | 0.361 |
|                           | 0.806   | 0.657        | 0.735       |                     | 1.672  | 318           |           | turn        | 0.361 |

1.588 318

1.999



0.324

depression

0.728

0.796

0.689

0.699

0.709

0.699

0.738

sens

0.537

0.610

0.707

0.561

0.610

0.683

0.683

0.758

0.774

0.625

0.667

0.746

0.787

# Case-Studies – General Populations

| 2<br>2<br>100<br>100<br>100<br>100<br>100<br>100<br>101<br>101<br>101 | 20005<br>110007<br>25780<br>12139<br>12188<br>12187<br>12663<br>12704<br>12652<br>12292<br>12291<br>2025<br>22672 | Ongoing characteristic<br>Ongoing characteristic<br>Brain MRI Acquis<br>Brain MRI Believ<br>Brain MRI Brain I<br>Brain MRI Reaso<br>Brain MRI Reaso<br>Brain MRI Reaso<br>Carotid ultrasound<br>Carotid ultrasound<br>Carotid ultrasound | s Email access<br>s Newsletter communications, d<br>ition protocol phase.<br>ed safe to perform brain MRI sca<br>MRI measurement completed<br>MRI measuring method<br>n believed unsafe to perform bra<br>n brain MRI not completed<br>n brain MRI not completed<br>Carotid ultrasound measuring<br>Carotid ultrasound measuring<br>Carotid ultrasound measuring<br>Carotid ultrasound measuris pack<br>Maximum carotid IMT (intima- | ate sent<br>n<br>in MRI<br>ent completed<br>method<br>sage<br>medial thickness | )#t120   | UK Biobank – disc<br>between HC, single<br>multiple comorbid<br>Predict likelihoods<br>developmental or a<br>disorders<br>Forecast cancer | riminate<br>e and<br>conditions<br>of various<br>aging |
|---|---|--|--|--|----------|---|--|
| 101   | 22675   | Carotid ultrasound   | Maximum carotid IMT (intima-   | medial thickness   | ) at 150 |   |  |
| degree<br>101<br>degree   | es<br>22678<br>es   | Carotid ultrasound   | Maximum carotid IMT (intima-   | Data<br>Source   | Sam      | nple Size/Data Type   | Summary  |
| 101   | 22681   | Carotid ultrasound   | Maximum carotid IMT (intima-   |  | _        |   |  |
| degre   | es  |  |  |  | Den      | nographics: > 500K cases  | The  |
| 101   | 22671   | Carotid ultrasound   | Mean carotid IMT (intima-med   |  | Clin     | ical data: > 4K features  | longitudinal   |
| 101   | 22674   | Carotid ultrasound   | Mean carotid IMT (intima-med   | LIK .  | Ima      | ging data: T1 rosting   | archive of   |
| 101   | 22677   | Carotid ultrasound   | Mean carotid IMT (intima-med   |  | iiiia    | ging uata. 11, resting-   |  |
| 101   | 22680   | Carotid ultrasound   | Mean carotid IMT (intima-med   | Biobank  | state    | e fMRI, task fMRI,  | the UK   |
| 101   | 22670   | Carotid ultrasound   | Minimum carotid IMIT (Intima-I   |  | T2       | FLAIR, dMRL SWI   | population   |
| degree  | 220222  | Constitution and   | Adiation and an actual IDAT (inting a  |  |          |   |  |
| 101<br>dogrou   | 22673   | Carolid ultrasound   | winimum caroud iwir (inuma-i   |  | Gen      |   | (110)  |
| 101   | -><br>  | Carotid ultracound   | Minimum carotid IMT (intima.)  | modial thicknose)  | at 210   |   |  |
| degree  | 22070   | Carolia all'asouna   |  | incular the kiness)  | 01210    | http://www.ukbioba  | ank.ac.uk  |
| 101   | 22679   | Carotid ultrasound   | Minimum carotid IMT (intima-   | medial thickness)  | at 240   | http://hd2k.org   |  |
| degree  | 25  |  |  | ,  |          | nttp://buzk.org   |  |
| 101   | 22682   | Carotid ultrasound   | Quality control indicator for IN   | IT at 120 degrees  |          |   |  |
| 101   | 22683   | Carotid ultrasound   | Quality control indicator for IN   | IT at 150 degrees  |          |   |  |
|   |   |  |  |  |          |   |  |

## Case-Studies – UK Biobank (Complexities)



## Case-Studies – UK Biobank – NI Biomarkers



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



# End-to-end Pipeline Workflow Solutions





| AHC AS   | SOCIATION O<br>EALTH CARE J  |   | )1u      | UNITED STATES D  | EPARTMENT OF LABOR   | ury tico).   | Dutu I  | <u>usion</u>              |
|--|--|---|----------|--|--|--|---|---------------------------|
| Search Resc<br>< last page<br>Advanced Search  | Census<br>ZIP Code<br>Statistics at  | People Busin<br>Us Dependent of Community<br>Census<br>Main COMMUNITY FACTS   | ess   c  | Home V Subjects V<br>Databases, V<br>TOP Picks   | OF LABOR STATIS  | Awzaw<br>STICS<br>Economic Releases - Studen<br>ors by Subject   | der   FAQE   Albert BLS   Carles to<br>Fridere to st   Laware to<br>Same RLS gov<br>RLS v Bets v<br>DAME con 1 1 1 10 | Address to Free Laders 64 |
| Home<br>About AHCJ<br>Advertise with AHCJ<br>Center for Excellence in<br>Health Care Journalism<br>Membership      | Census 2010     Statistics at <u>2007 Econom</u> FactFinder. 20     2002 Econom  | Search - Use the optio<br>Your Selections<br>Search using<br>Clear al selections and<br>start a new search  | ns on ti | PUBLIC DATA APT<br>DISCONTINUED DATABASES<br>PAQS<br>SPECIAL NOTICES<br>MORE SOURCES OF DATA | on This Page:         Inflation & Prices         Employment         Unemployment         Vinemployment         Pay & Benefits         Spending & Time Use              | <ul> <li>Productivity</li> <li>Workplace. Injuries</li> <li>International</li> <li>Employment Projections</li> <li>Regional Resources</li> </ul> | <ul> <li>Historical News<br/>Release Tables</li> <li>Maps</li> <li>Calculators</li> <li>Public Data API</li> </ul>    |                           |
| Training<br>Resources<br>Core Topic: Health &<br>Core Topic: Aging<br>Core Topic: Oral Hea<br>Core Topic: Med. Stu | FactFinder an<br>Census.<br><u>1997 Econom</u><br>industry by en<br>18 sectors co<br>1992 Econom<br>they are not e<br><u>ZIP Business</u><br>the Economic<br>available in at   | Search using the options below:<br>Topics<br>(dep. income, year, dataset)<br>Geographies<br>(dates, countes, piaces,)<br>Race and Ethnic Groups<br>(ace, ancesity, the) |          |  | Inflation & Prices<br>Database Na<br>Prices - Consumer   | ime Special<br>Notice  | Top One Multi-<br>Picks Screen Screen   | Tables                    |
|  | Classified by 5<br>Classified by | Industry Codes<br>(NACS industry,)<br>EEO Occupation Codes<br>(executives, analysis,)   |          |  | (Consumer Pinor Index - CPI)<br>Urban Wage Earners and Clerical II<br>(Consumer Pinor Index - CPI)<br>All Urban Consumers (Chained CPI<br>(Consumer Pinor Index - CPI) | Norkers (Current Series)   |   |                           |
| http://socr.   | umich.ed   | u/HTML5/Da  | shbo     | oard   | Average Price Data<br>Second Later Data<br>(Connect Later Data - CA)<br>VIL OPER Commence (Convert CA)   | 0  |   |                           |

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

## Compressive Big Data Analytics (CBDA)

- Foundation 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 nonparametric 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)



## FAIR Data & Open-Science Principles

Open Science

Finge

sable

Access

Interop

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



## 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





## **Clinical Decision Support**

Personalized medicine – Traumatic Brain Injury (TBI)



## Acknowledgments 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/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 STATISTICS



- Complex DB Search, retrieval (IDA)
- Multidimensional data visualization (MotionCharts, BrainViewer, R)
- o 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)



