Compressive Big Data 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 Biomedical Data
- Data science & predictive analytics
- □ Compressive Big Data Analytics (CBDA)
- Case-studies
 - Applications to Neurodegenerative Disease
 - Data Dashboarding





Driving Biomedical/Health Challenges

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.



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 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)





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</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 – Parkinson's Disease (1)

- 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), in press





Case-Studies – Parkinson's Disease (1)

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), in press



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 – UK Biobank (Complexities)



Case-Studies – UK Biobank – NI Biomarkers









Case	-Stu	dies -	- UK Biobank –	Results	5
Yurishide Sex Fenale Masie Sensitivity/hurt feelings Yes No Worrite/anolous feelings	Cluster 1 1,134 (24.7%) 3,461 (75.3%) 2,142 (47.9%) 2,332 (52.1%)	4,062 (76. 5) 1,257 (23. 5) 3,023 (58. 5) 2,151 (41. 5)			
res No Risk taking Yes	2,173 (48.2%) 2,337 (51.8%) 1,378 (31.0%)	2,995 (57. 5) 2,208 (42. 5) 1,154 (22. 5)	Variable	Cluster 1	Cluster 2
NO Guilty feelings Yes No Seen doctor for nerves, anxiety, tension or depression	3,064 (69.0%) 1,100 (24.4%) 3,417 (75.6%)	1,697 (32. i) 3,536 (67. i)	Sex Female	1,134 (24.7%)	4,062 (76.4%)
Yes No Alcohol usually taken with meals Yes	1,341 (29.3%) 3,237 (70.7%) 1,854 (66.7%)	1,985 (37. 5) 3,310 (62. 5) 2,519 (76. 5)	Male	3,461 (75.3%)	1,257 (23.6%)
No Snoring Yes No	924 (33.3%) 1,796 (41.1%) 2,577 (58.9%)	771 (23.45 1,652 (33.5) 3,306 (66.5)	•••		_
Worry too long after embarrassment Yes No Miserableness	1,978 (44.3%) 2,491 (55.7%)	2,675 (52. 5) 2,462 (47. 5)	Yes	751 (16.6%)	1,071 (20.8%)
Yes No Ever highly irritable/argumentative for 2 days Yes No	1,715 (37.7%) 2,829 (62.3%) 485 (10.7%) 4.038 (89.3%)	2,365 (45.5) 2,882 (54.5) 749 (14.5%) 4,418 (85.5)		5,705 (85.4%)	4,070 (75.276)
Nervous feelings Yes No Ever depressed for a whole week	751 (16.6%) 3,763 (83.4%)	1,071 (20. 5) 4,076 (79. 5)	Frequency of tiredness/lethargy in		
Yes No Ever unenthusiastic/disinterested for a whole week Yes	2,176 (48.1%) 2,347 (51.9%) 1,346 (30.3%)	2,739 (52. 5) 2,438 (47. 5) 1,743 (34. 5)	last 2 weeks Not at all	2,402 (53.0%) 1,770 (39.0%)	2,489 (47.8%) 2,127 (40.9%)
No Sleepless/insomnia Never/rarely Sometimes	3,089 (69.7%) 1,367 (29.8%) 2,202 (47.9%)	3,344 (65. 5) 1,181 (22. 5) 2,571 (48. 5)	Several days More than half the days	187 (4.1%1) 177 (3.9%)	300 (5.8%) 287 (5.5%)
Usually Getting up in morning Not at all easy Not very easy Fold easy	1,024 (22.3%) 139 (3.1%) 538 (11.9%) 2 227 (51.4%)	1,563 (29. 5) 249 (4.7% 830 (15.81 2.662 (50. 1)	Nearly everyday Alcohol drinker status		
Very easy Nap during day Never/rarely Sometimes	2,497 (54.5%) 2,497 (54.5%) 1,774 (38.8%)	1,505 (38. 5) 3,238 (61. 5) 1,798 (34. 5)	Never	81 (1.8%)	179 (3.4%)
Usually Frequency of tiredness/lethargy in last 2 weeks Not at all Several days	2,402 (53.0%) 1,770 (39.0%)	228 (4.3% 2,489 (47. 5) 2,127 (40. 5)	Current	4,429 (96.4%)	4,992 (93.9%)
More than half the days Nearly everyday Alcohol drinker status Never	187 (4.1%1) 177 (3.9%) 81 (1.8%)	300 (5.8% 287 (5.5% 179 (3.4%			
Previous Current	83 (1.8%) 4,429 (96.4%)	146 (2.7%)			



Decision tree illustrating a simple clinical decision support system providing machine guidance for identifying <u>depression feelings</u> based on categorical variables and neuroimaging biomarkers. In each terminal node, the y vector includes the percentage of subjects being labeled as "no" and "yes", in this case, answering the question "Ever depressed for a whole week." The p-values listed at branching nodes indicate the significance of the corresponding splitting criterion.

Case-Studies	– UK	Biobank –	Resul	ts
	Accuracy	95% CI (Accuracy)	Sensitivity	Specificity
Sensitivity/hurt feelings	0.700	(0.676, 0.724)	0.657	0.740
Ever depressed for a whole week	0.782	(0.760, 0.803)	0.938	0.618
Worrier/anxious feelings	0.730	(0.706, 0.753)	0.721	0.739
Miserableness	0.739	(0.715, 0.762)	0.863	0.548

Cross-validated (random forest) prediction results for four types of mental disorders



End-to-end Pipeline Workflow Solutions



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 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

Dinov, 2016, PMID: 26998309;

Marino, et al., (pending



Synergies with Compressive Sensing

• Define the nested sets

subsequent steps)

 $S_k = \{x: \|x\|_o \stackrel{\text{def}}{=} |supp(x)| \le k\},\$

where the data x, as a vector or tensor, has at most k non-trivial elements. Note that if $x, z \in S_k$, then $x + z \in S_{2k} \supseteq S_k$

• If $\Phi_{n \times n} = (\varphi_1, \varphi_2, \varphi_3, ..., \varphi_n)$ represents an orthonormal basis, the data may be expressed as $x = \Phi c$, where $c_i = \langle x, \varphi_i \rangle$, i.e., $c = \Phi^T x$, and $||c||_o \le k$. Even if x is not strictly sparse, its representation c may be sparse. For each dataset, we can assess and quantify the error of approximating x by an optimal estimate $\hat{x} \in S_k$ by computing

$$\sigma_k(x)_p = \min_{\hat{x} \in S_k} \|x - \hat{x}\|_p$$



Synergies with Compressive Sensing

• In compressive sensing, if $x \in \mathbb{R}^n$, and we have a data stream generating m linear measurements, we can represent y = Ax, where $A_{m \times n}$ is a dimensionality reducing matrix ($m \ll n$), i.e.,

$$A_{m \times n}: \mathbb{R}^n \to \mathbb{R}^m$$

• The null space of A, $N(A) = \{z \in R^n : Az = 0 \in R^m\}$

A uniquely represents all $x \in S_k \Leftrightarrow N(A)$ contains no vectors in S_{2k} .



• The spark of a matrix A represents the smallest number of columns of A that are linearly dependent. If $A_{m \times n}$ is a random matrix whose entries are independent and identically distributed, then spark(A) = m + 1, with probability 1.





Synergies with Compressive Sensing

◦ If the entries of *A* are chosen according to a sub-Gaussian distribution, then with high probability, for each *k*, there exists $\delta_{2k} \in (0,1)$ such that

 $(1 - \delta_{2k}) \|x\|_2^2 \le \|Ax\|_2^2 \le (1 + \delta_{2k}) \|x\|_2^2$ (1) for all $x \in S_{2k}$ (RIP=Restricted isometry property)

• When we know that the original signal is sparse, to reconstruct *x* given the observed measurements *y*, we can solve the optimization problem:

$$\hat{x} = \arg\min_{z:Az=v} \|z\|_0$$

Synergies with Compressive Sensing

- Linear programming may be used to solve the optimization problem if we replace the zero-norm by its more tractable convex approximation, the l_1 -norm, $\hat{x} = \arg\min_{x \in I = 1} ||z||_1$
- Given that $A_{m \times n}$ has the above property and $\delta_{2k} < \sqrt{2} 1$, if we observe y = Ax, then the solution \hat{x} satisfies

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

• Thus, in compressive sensing applications, if $x \in S_k$ and A satisfies the RIP, condition (1), we can recover any k-sparse signal x exactly (as $\sigma_k(x)_1 = 0$)

using only $O(k \log(n/k))$ observations, since $m = O\left(\frac{k \log(n/k)}{\delta_{2k}^2}\right)$

• Finally, if $A_{m \times n}$ is random (e.g., chosen according to a Gaussian distribution) and $\Phi_{n \times n}$ is an orthonormal basis, then $A_{m \times n} \times \Phi_{n \times n}$ will also have a Gaussian distribution, and if *m* is large, $A' = A \times \Phi$ will also satisfy the RIP condition (1) with high probability.



Compressive Big Data Analytics (CBDA)

- □ Start with a generic dataset: $[X, Y, X_{val}, Y_{val}]$, let $C^j = [X^j, Y^j, X^j_{val}]$, j = 1, ..., M.
- □ Define machine learning predictor as ML: $ML(C^j)$: $R^{n_j \times k_j} \times R^{n_j} \times R^{m \times k_j} \to R^m$, $ML([X^j, Y^j, X^j_{val}]) = \hat{Y^j}$.
- □ Define performance metric as: $\tau(ML(C^j), Y_{val}): \mathbb{R}^m \times \mathbb{R}^m \to \mathbb{R}, \tau(Y^j, Y_{val}) = c_j, j = 1, 2, ..., M.$
- **Rank the samples**: $\{c_{(j)}\} = O^q(\{c_j\}), j = 1, 2, ..., q$.
- □ Set the feature values: $s_j^{(i)} = Dirichlet_j(f_i) = \begin{cases} 0, & f_i \notin sample \ j \\ 1, & f_i \in sample \ j \end{cases}$
- □ Set a metric as: $F \in R^{q \times K}$, $\forall b_{ii} \in F$, $b_{ii} = s_i^{(i)}$.
- **Count the occurrence**: $S_i = \sum_{i=1}^{q} b_{ii}$, i = 1, 2, ..., K.
- **•** Feature mining: $S_{(i)} = O^K(S_i), i = 1, 2, ..., K, \Omega^* = \{f_1, ..., f_i, ..., f_K\}.$
- □ Inference: let $C^p = [\phi_p X, \phi_p Y, \phi_p X_{val}]$ and $[ML_p(C^p): R^{n \times k_p *} \times R^{n \times 1} \times R^{m \times k_{p^*}} \rightarrow R^{m \times 1}, ML([\phi_p X, \phi_p Y, \phi_p X_{val}]) = Y_{val}^p, p = 5, 10, \dots, \hat{p}, \tau(Y_{val}^p, Y_{val}): R^m \times R^m \rightarrow R, \tau(ML([\phi_p X, \phi_p Y, \phi_p X_{val}]), Y_{val}) = best\tau(ML([\phi_p X, \phi_p Y, \phi_p X_{val}]), Y_{val}).$
- \Box Then Φ_{p^*} is the final dictionary we need.
- □ Performance assessment: Once we obtain ϕ_{p^*} , and have X_{exp} which is under prediction. Then perform ϕ_{p^*} and SuperLearner algorithm: $Y_{exp} = ML([\phi_{p^*}X, \phi_{p^*}Y, \phi_{p^*}X_{exp}])$, where Y_{exp} represents the expected results.





CBDA Results: Simulated Data

Knockoff of Null (left) vs Binomial (right) Data

Panels A, C and E show the correspondent histograms generated from the Knockoff Filter algorithm on the three Null datasets.

Panels B, D and F show the correspondent histograms generated from the Knockoff Filter algorithm on the three Binomial datasets.

Performance metric: MSE





CBDA Results: Simulated Data

CBDA Results: Null (left) vs Binomial (right) Data

Panels A, C and E show the correspondent histograms of the CBDA Results on the three Null datasets.

Panels B, D and F show the correspondent histograms of the CBDA results on the three Binomial datasets.

Performance metric: MSE



CBDA Results: Biomed Data (ADNI)

	Reference			
Prediction	AD	MCI	Normal	Ê
AD	69	17	1	Õ
MCI	12	243	8	
Normal	0	9	140	nu
Ov	erall Statistics	;		ltin
Accuracy	0.9058 [95%	6 CI = (0.8767,	0.93)]	P
No Information Rate	0.5391			nia
P-Value [Acc > NIR]	<2e-16			0
Карра	0.8426			lag
McNemar's Test P-Value	0.589			šit
Statistics	by Diagnosti	c Class		lic
	AD	MCI	Normal	atic
Sensitivity	0.8519	0.903	3 0.9396	Ĕ
Specificity	0.9569	0.913	0.9743	ſe
Positive Pred Value	0.7931	0.924	0.9396	sul
Negative Pred Value	0.9709	0.8898	8 0.9743	ts
Prevalence	0.1623	0.539	0.2986	$\widehat{}$
Detection Rate	0.1383	0.4870	0.2806	9
Detection Prevalence	0.1743	0.527	0.2986	É
Balanced Accuracy	0.9044	0.9082	2 0.9569	

Acknowledgments

Funding

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

http://SOCR.umich.edu

The Elsie Andresen Fiske Research Fund

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

- LONVINI: 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: Cathle Spino, Chuck Burant, Ben Hampstead, Stephen Goutman, Stephen Strobbe, Hiroko Dodge, Hank Paulson, Bill Dauer, Brian Athey

