

DataSifter: Sharing of Sensitive Information via Statistical Obfuscation

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>

Slides Online:
"SOCR News"



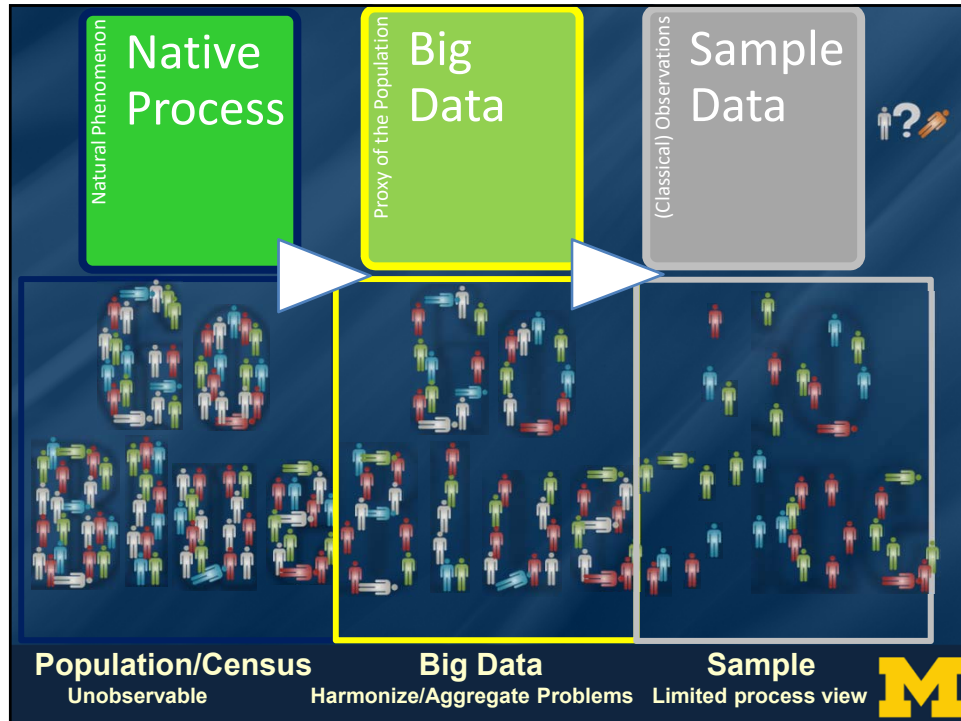
SCHOOL OF MEDICINE
UNIVERSITY OF MICHIGAN

STATISTICS ONLINE COMPUTATIONAL RESOURCE (SOCR)

Outline

- ☐ Driving biomedical & health challenges
- ☐ Common characteristics of Big Neuroscience Data
- ☐ ϵ -Differential privacy (DP); Homomorphic encryption
- ☐ *DataSifter: Statistical obfuscation*
- ☐ Case-studies
 - ☐ Applications to Neurodegenerative Disease (Udall/MADC)
 - ☐ Autism Brain Imaging Data Exchange (ABIDE)
 - ☐ Population Census-like Neuroscience





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
Time	Techniques accounting for longitudinal patterns in the data
Incomplete	Reliable management of missing data

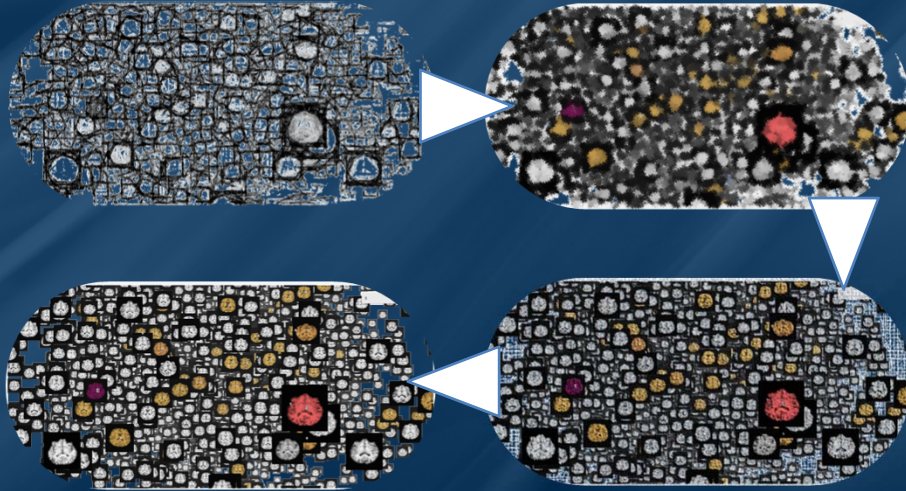
Example: analyzing observational data of 1,000's Parkinson's disease patients based on 10,000's signature biomarkers derived from multi-source imaging, genetics, clinical, physiologic, phenomics and demographic data elements

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



Multiscale/Multimodal NI Data



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



ϵ -Differential privacy (ϵ DP) vs. Homomorphic encryption (HE)

Category	ϵ DP	HE
Goal	Mine information in a DB without compromising privacy; no access to inspect individual DB entries	Provide a secure encryption allowing program execution on encrypted data; encrypt results, interpretation requires ability to decrypt the data



ϵ -Differential privacy (ϵ DP)

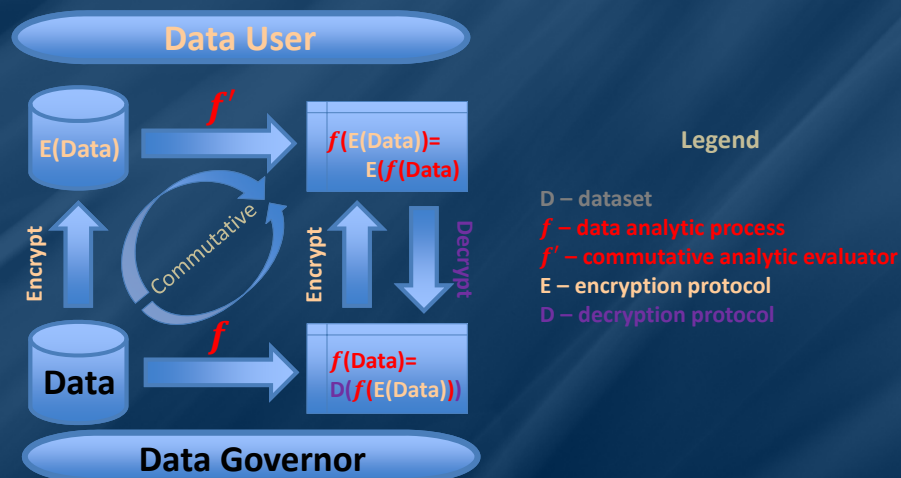
- ❑ **Data-features:** $\{C_1, C_2, \dots, C_k\}$, categorical or numerical.
- ❑ **DB** = list of cases $\{x_1, x_2, \dots, x_n\}$, $x_i \in C_1 \times C_2 \times \dots \times C_k$, $1 \leq i \leq n$.
- ❑ ϵ -Differential privacy relies on adding noise to data to protect the identities of individual records. An **algorithm** f is ϵ -differentially private if for all possible inputs (datasets/DBs) D_1, D_2 that differ on a single record, and all possible f outputs, y , the probability of correctly guessing D_1 knowing y is not significantly different from that of D_2 :

$$\frac{P(f(D_1) = y)}{P(f(D_2) = y)} \leq e^\epsilon, \quad \forall y \in \text{Range}(f).$$
- ❑ The global sensitivity of f is the smallest number $S(f)$, such that $\forall D_1, D_2$ that differ on at most one element $\|f(D_1) - f(D_2)\|_1 \leq S(f)$
- ❑ There are many differentially private algorithms, e.g., random forests, decision trees, k-means clustering, etc.
- ❑ E.g., $f: D = \text{DB} \rightarrow R^m$, the algorithm outputting $f(D) + (y_1, y_2, \dots, y_m)$, with $y_i \in \text{Laplace}\left(\mu = 0, \sigma = \sqrt{2} \frac{S(f)}{\epsilon}\right)$, $\forall i$ is ϵ -differentially private

Dwork, LNCS, 2008



Homomorphic encryption (HE)



Rivest & Adleman, Academic Press, 1978



DataSifter

- ❑ DataSifter is an iterative statistical computing approach that provides the data-governors controlled manipulation of the trade-off between sensitive information obfuscation and preservation of the joint distribution.
- ❑ The DataSifter is designed to satisfy data requests from pilot study investigators focused on specific target populations.
- ❑ Iteratively, the DataSifter stochastically identifies candidate entries, cases as well as features, and subsequently selects, nullifies, and imputes the chosen elements. This statistical-obfuscation process relies heavily on nonparametric multivariate imputation to preserve the information content of the complex data.

<http://DataSifter.org>

US patent #16/051,881

Marino, Zhou, *et al.*, in review (2018)



DataSifter

- ❑ dataSifter() R method implementation and detailed description are available on our GitHub repository (<https://github.com/SOCR/DataSifter>).
- ❑ Data-sifting different data archives requires customized parameter management. Five specific parameters mediate the balance between protection of sensitive information and signal energy preservation.

Obfuscation level	k_0	k_1	k_2	k_3	k_4
None	0	0	0	0	0
Small	0	0.05	1	0.1	0.01
Medium	1	0.25	2	0.6	0.05
Large	1	0.4	5	0.8	0.2
Indep	Output synthetic data with independent features				

k_0 : A Boolean; obfuscate the unstructured features?

k_1 : proportion of artificial missing data values that should be introduced

k_2 : The number of times to iterate

k_3 : The fraction of structured features to be obfuscated in all the cases

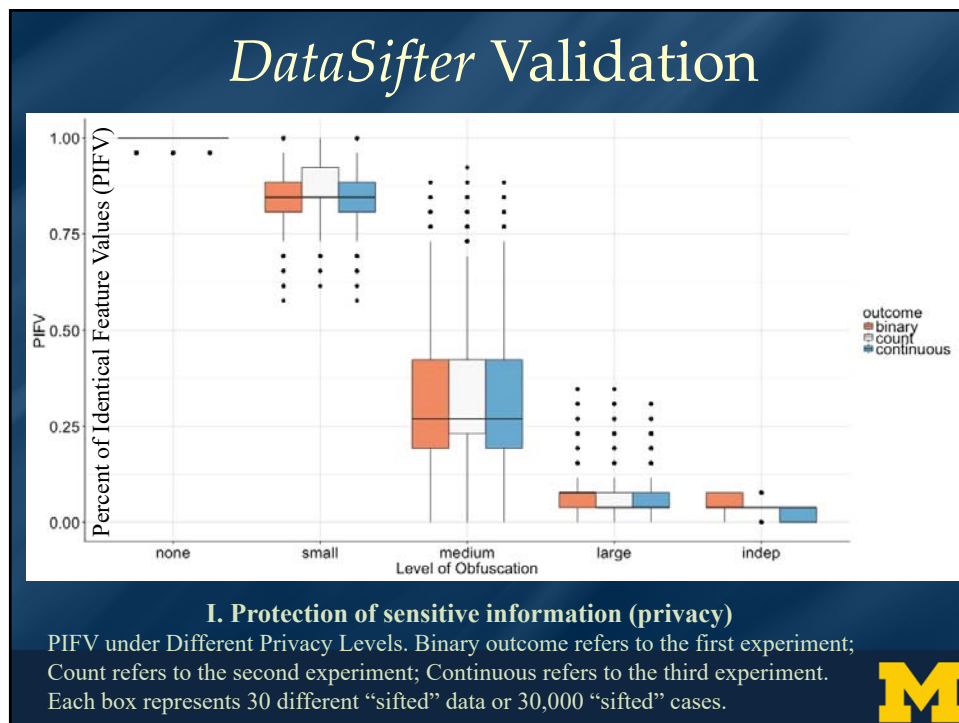
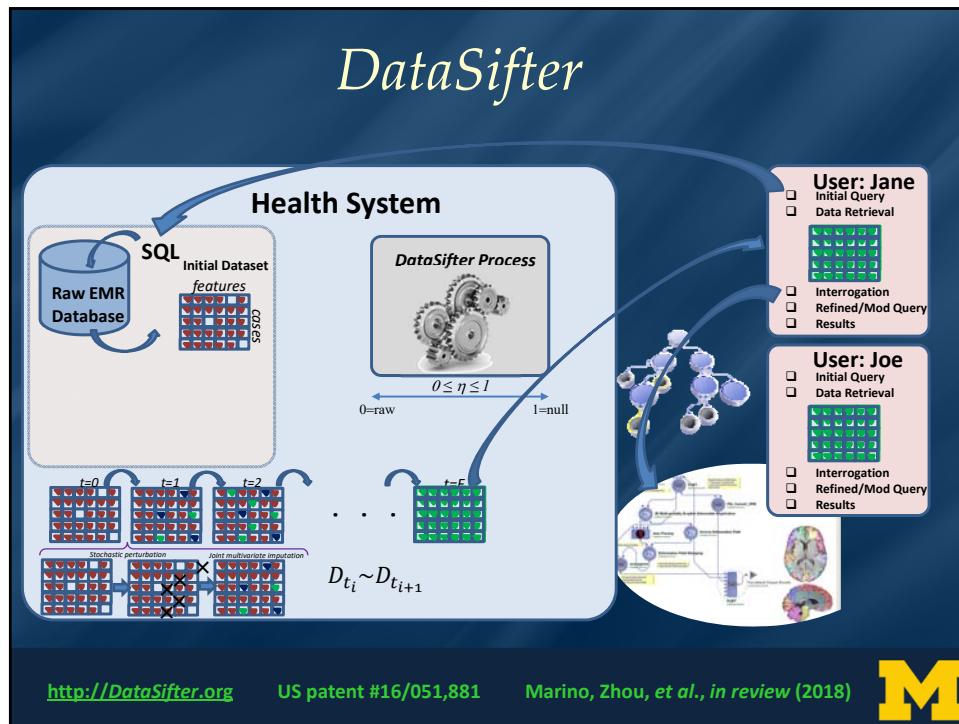
k_4 : The fraction of closest subjects to be considered as neighbours of a given subject

<http://DataSifter.org>

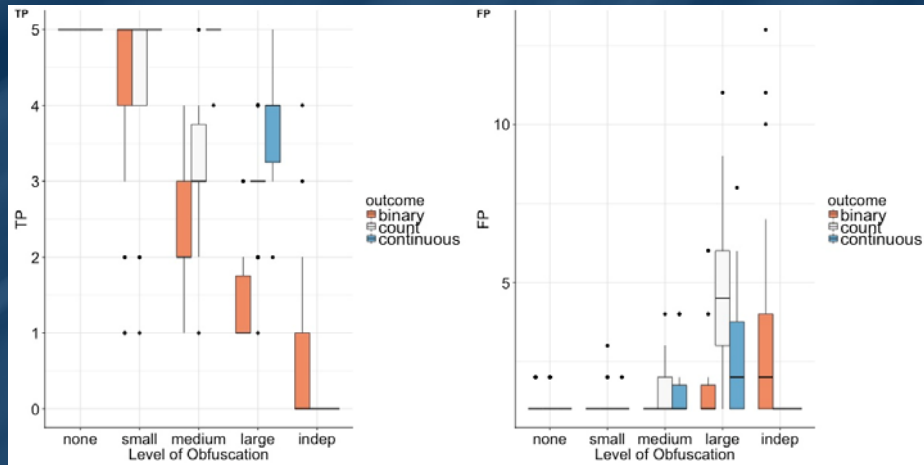
US patent #16/051,881

Marino, Zhou, *et al.*, in review (2018)





DataSifter Validation



II. Preserving utility information of the original dataset

Logistic Model with Elastic Net Signal Capturing Ability. TP is the number of true signals (total true predictors = 5) captured by the model. FP is the number of null signals that the model has falsely selected (total null signals=20).



DataSifter Validation

III. Clinical Data Application: Using DataSifter to Obfuscate the ABIDE Data

Comparing the Original and “Sifted” Data for the 22nd ABIDE Subject

	Output	Sex	Age	Acquisition Plane	IQ	thick_std_ct x .lh.cuneus	curv_ind_ctx _lh_G_front_inf.Triangul	gaus_curv_ctx.lh.medialorbitofrontal	curv_ind_ctx _lh_S_interm_prim.Jensen
original	Autism	M	31.7	Sagittal	131	0.475	2.1	0.315	NA
none	Autism	M	31.7	Sagittal	131	0.475	2.1	0.315	0.51
small	Autism	M	31.7	Sagittal	131	0.475	2.1	0.315	0.4589
medium	Autism	M	31.7	Sagittal	111	0.548	2.85	0.315	0.463
large	Control	M	18.2	Sagittal	104	0.5347	3.198	0.1625	0.4524
indep	Control	M	15.4	Coronal	104	0.4842	3.383	0.1079	1.002



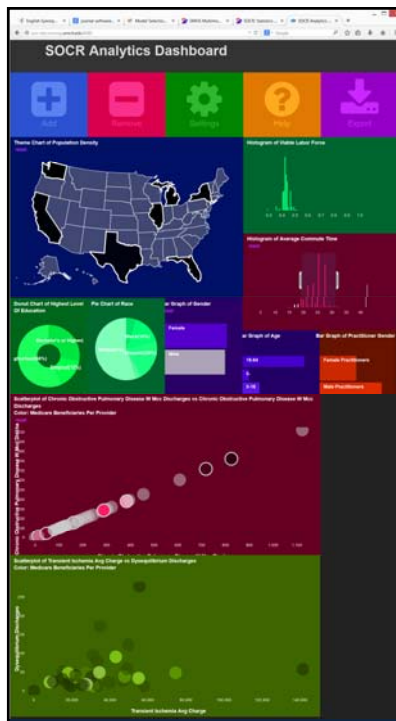
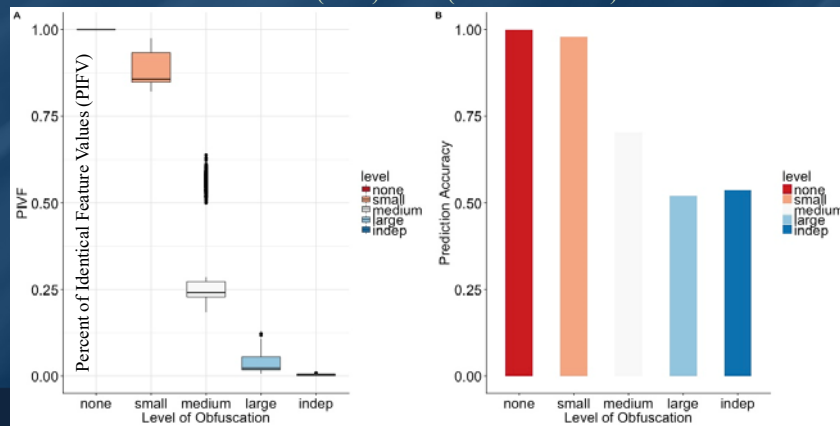
DataSifter Validation

IV. Clinical Data Application: Using DataSifter to Obfuscate the ABIDE Data

PIFVs for ABIDE under different levels of DataSifter obfuscations. Each box represents

1098 subjects among the ABIDE sub-cohort

Random forest prediction of binary clinical outcome - autism spectrum disorder (ASD) status (ASD vs. control)



SOCR Big Data Dashboard

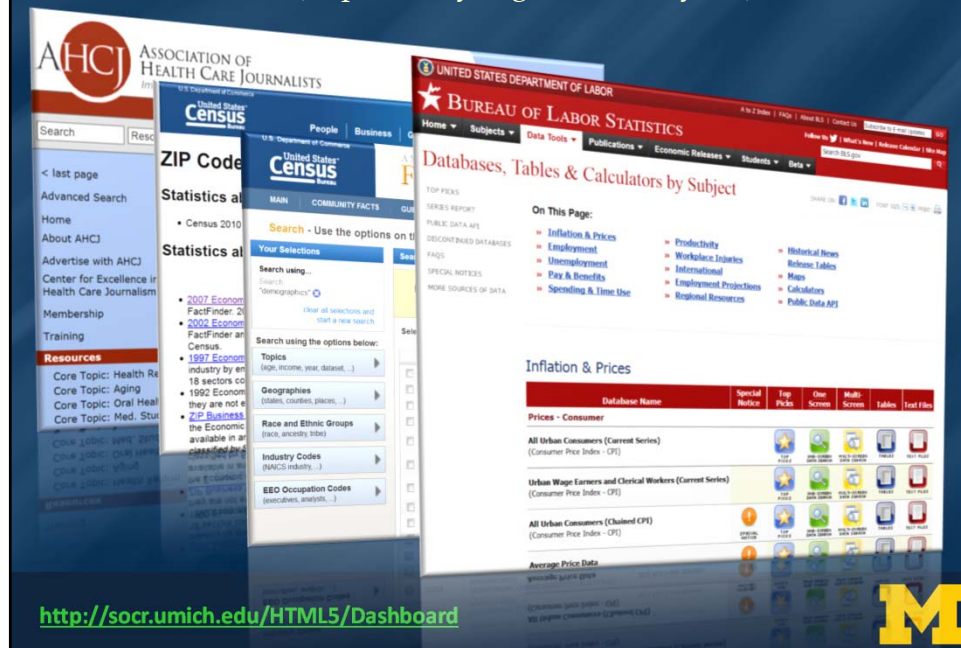
<http://socr.umich.edu/HTML5/Dashboard>

- ☐ 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

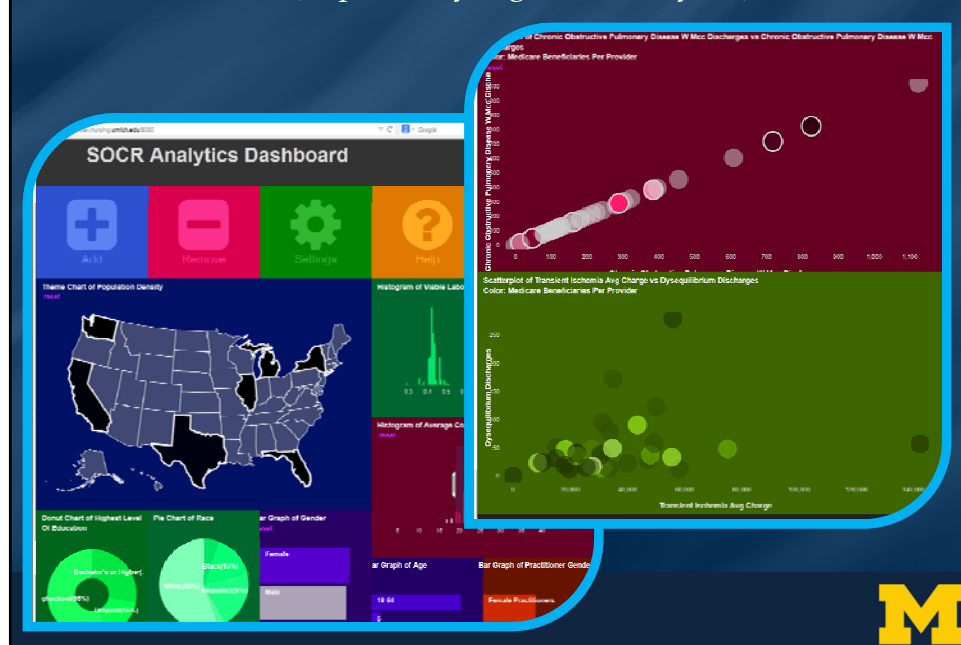
Husain, et al., 2015, PMID:26236573



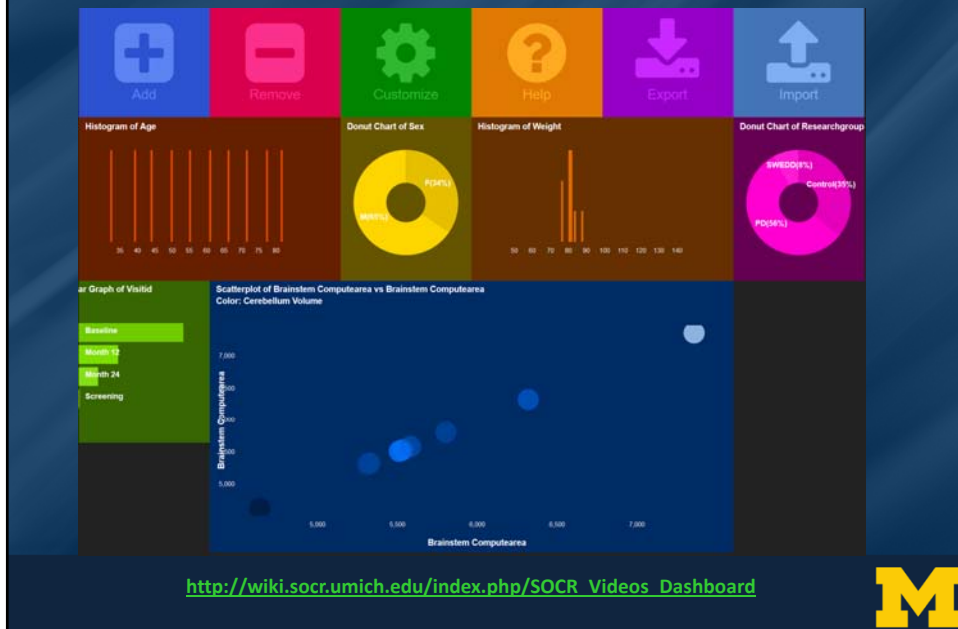
SOCR Dashboard (Exploratory Big Data Analytics): Data Fusion



SOCR Dashboard (Exploratory Big Data Analytics): Associations



SOCR Dashboard (Exploratory Big Data Analytics): Udall PD Data



Data Science & Predictive Analytics

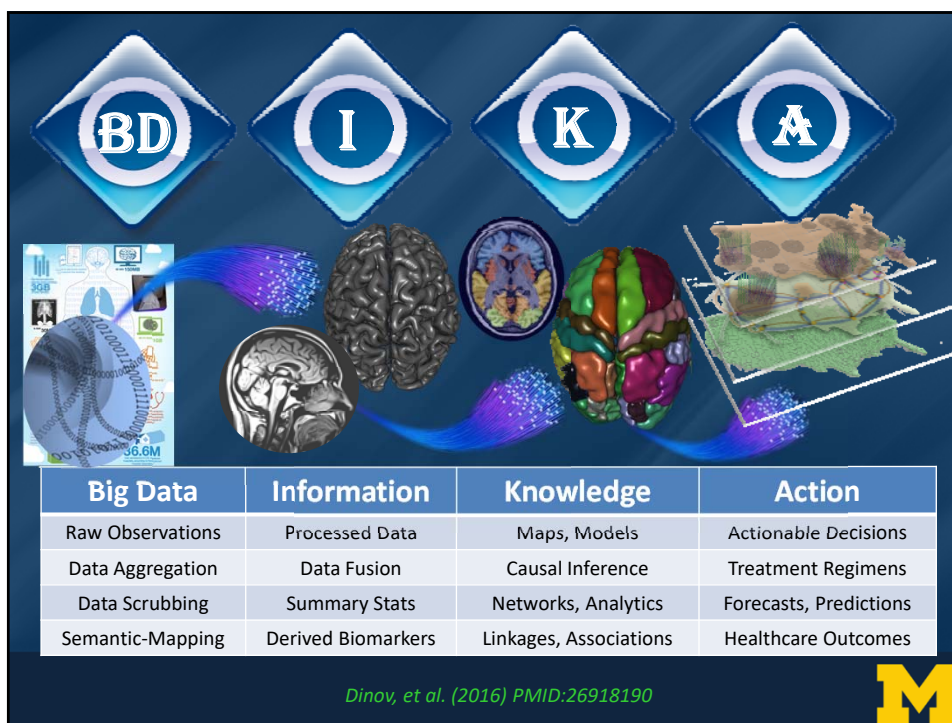
- ❑ **Data Science:** an emerging extremely transdisciplinary field - bridging between the theoretical, computational, experimental, and applied areas. Deals with enormous amounts of complex, incongruent and dynamic data from multiple sources. Aims to develop algorithms, methods, tools, and services 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

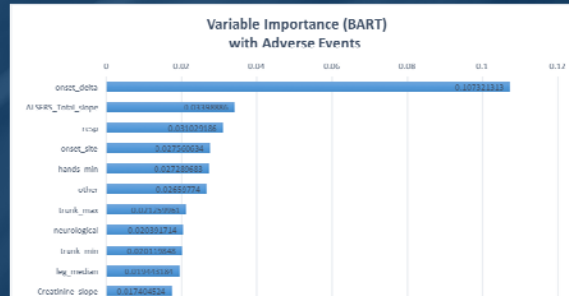
- ❑ Identify predictive classifiers to detect, track and prognosticate the progression of ALS (in terms of clinical outcomes like ALSFRS and muscle function)
- ❑ Provide a decision tree prediction of adverse events based on subject phenotype and 0-3 month clinical assessment changes

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 data</u> : <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)

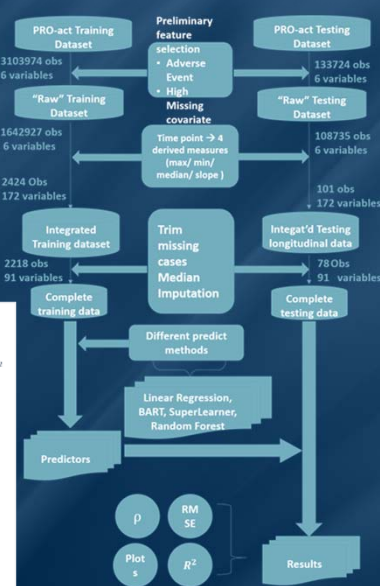
Tang, et al. (2018), in review

Case-Studies – ALS

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

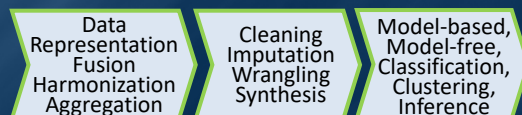


Methods	Linear Regression	Random Forest	BART	SuperLearner
R-squared	0.081	0.174	0.225	0.178
RMSE	0.619	0.587	0.568	0.585
Correlation	0.298	0.434	0.485	0.447



Case-Studies – ALS

- **Main Finding:** predicting univariate clinical outcomes may be challenging, the (information energy) signal is very weak. We can cluster ALS patients and generate evidence-based ALS hypotheses about the complex interactions of *multivariate factors*
- **Classification vs. Clustering:**
 - Classifying univariate clinical outcomes using the PRO-ACT data yields only marginal accuracy (about 70%).
 - Unsupervised clustering into sub-groups generates stable, reliable and consistent computable phenotypes whose explication requires interpretation of multivariate sets of features

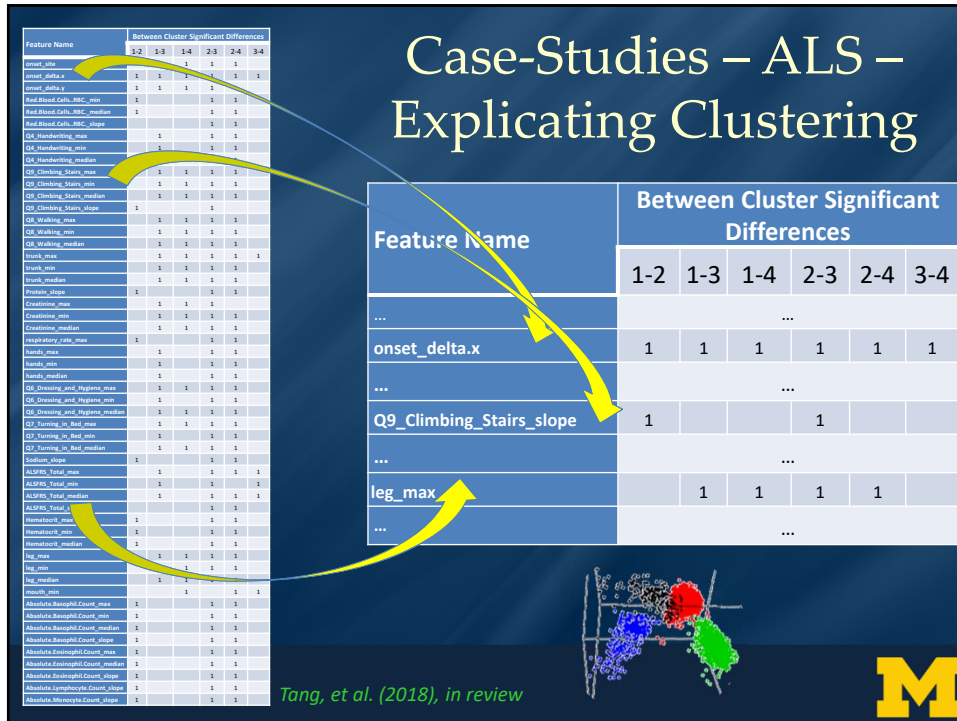


Cluster	Consistency	Variance	Cluster-Size	Silhouette
1	1	0	565	0.58
2	0.986	0.018	427	0.63
3	0.956	0.053	699	0.5
4	0.985	0.018	733	0.5

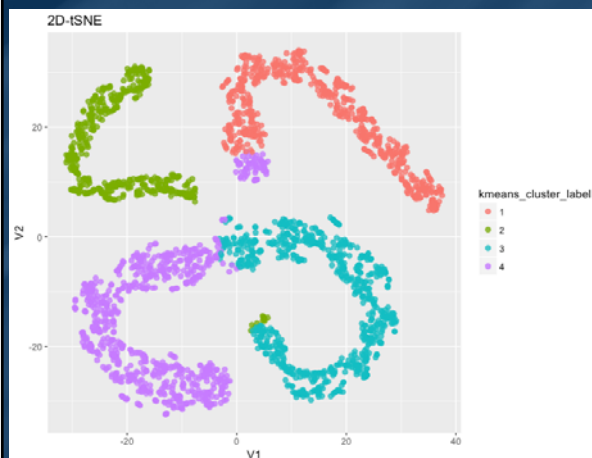
Tang, et al. (2018), in review



Case-Studies – ALS – Explicating Clustering



Case-Studies – ALS – Dimensionality Reduction



2D t-SNE Manifold embedding

Learn a mapping: $f: R^n \xrightarrow{n \gg d} R^d$
 $\{x_1, x_2, \dots, x_n\} \rightarrow \{y_1, y_2, \dots, y_d\}$
 preserves closely the *original distances*, $p_{i,j}$ and represents the *derived similarities*, $q_{i,j}$ between pairs of embedded points:

$$q_{i,j} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|y_i - y_k\|^2)^{-1}}$$

$$\min_f KL(P||Q) = \sum_{i \neq j} p_{i,j} \log \frac{p_{i,j}}{q_{i,j}}$$

$$0 = \frac{\partial KL(P||Q)}{\partial y_i} = 2 \sum_j (p_{i,j} - q_{i,j}) f(|x_i - x_j|) u_{i,j}$$

$$f(z) = \frac{z}{1+z^2} \text{ and } u_{i,j} \text{ is a unit vector from } y_j \text{ to } y_i.$$

Tang, et al. (2018), in review

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 forecast 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)



Case-Studies – Parkinson's Disease



Falls in PD are extremely difficult to predict ...



Case-Studies – Parkinson's Disease

Method	acc	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)



Open-Science & Collaborative Validation

End-to-end Big Data analytic protocol jointly processing complex imaging, genetics, clinical, demo data for assessing PD risk

- 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

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

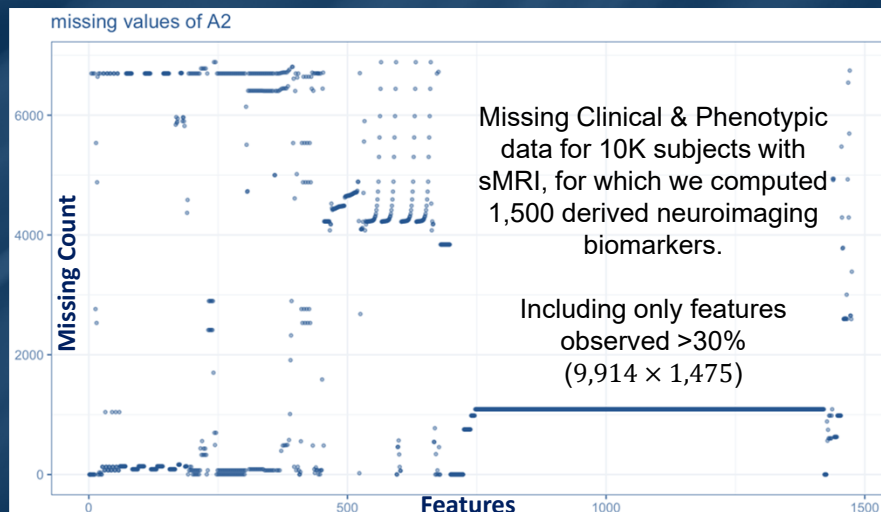
Data Source	Sample Size/Data Type	Summary
UK Biobank	Demographics: > 500K cases Clinical data: > 4K features Imaging data: T1, resting-state fMRI, task fMRI, T2_FLAIR, dMRI, SWI Genetics data	The longitudinal archive of the UK population (NHS)

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

<http://bd2k.org>



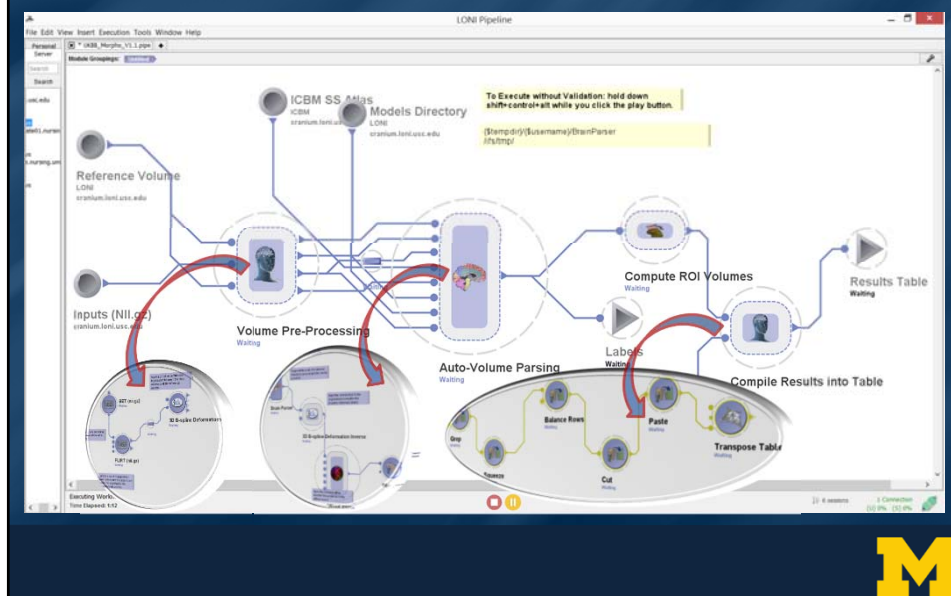
Case-Studies – UK Biobank (Complexities)



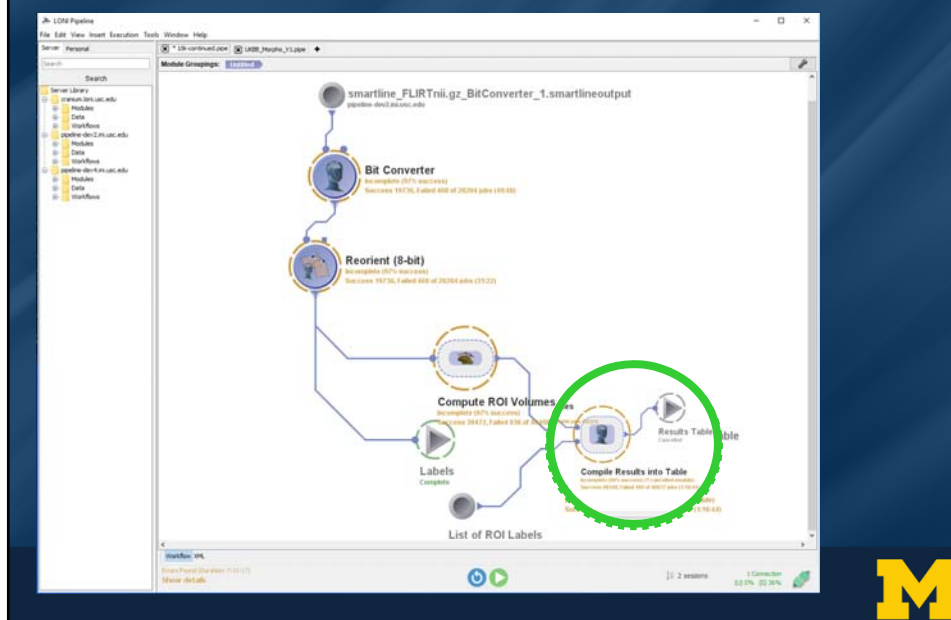
Zhou, et al. (2018), in review

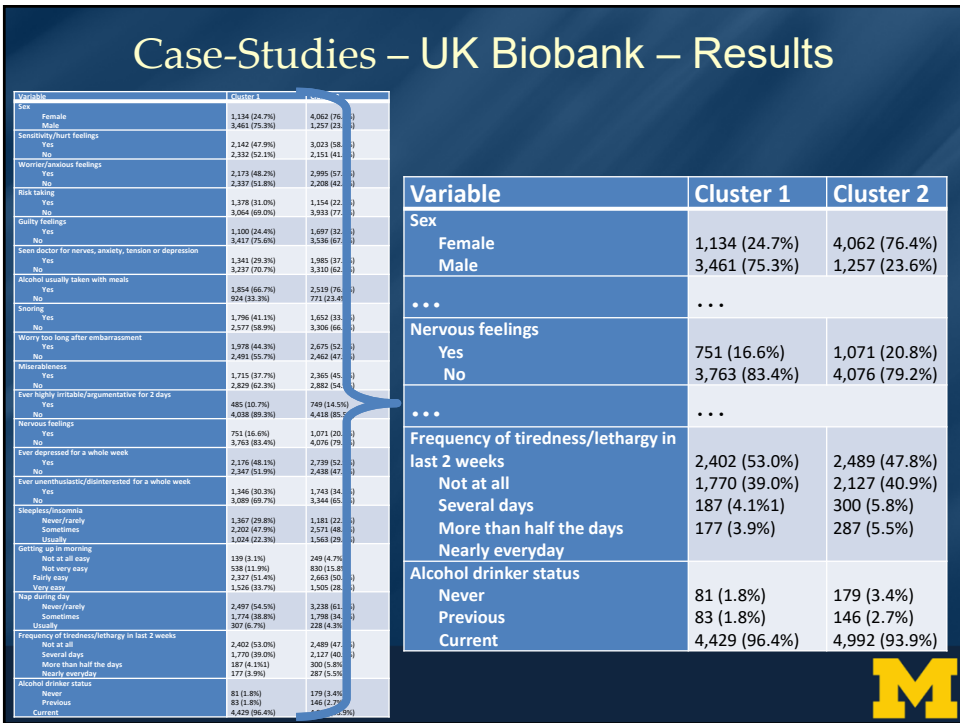
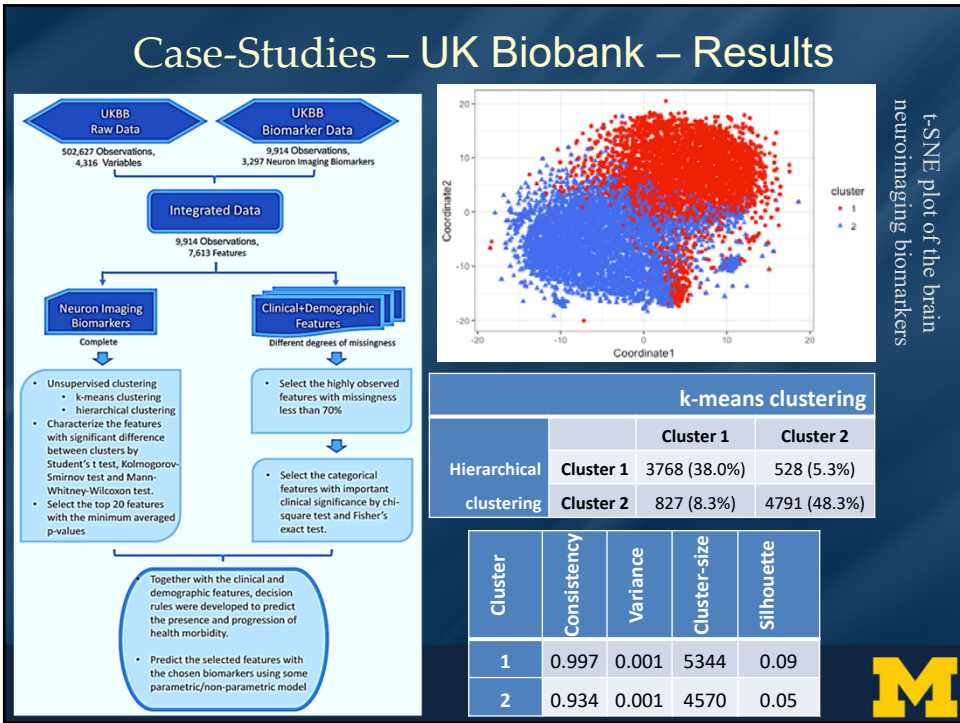


Case-Studies – UK Biobank – NI Biomarkers

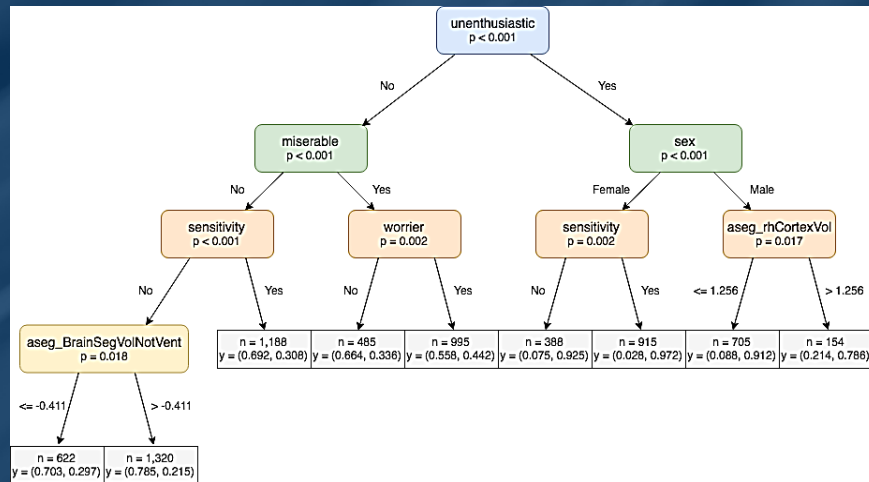


Case-Studies – UK Biobank – Successes/Failures





Case-Studies – UK Biobank – Results



Decision tree illustrating a simple clinical decision support system providing machine guidance for identifying **depression feelings** 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 – Results

	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

Zhou, et al. (2018), in review



Acknowledgments

Slides Online:
"SOCR News"

US patent #16/051,881

Funding

NIH: P20 NR015331, U54 EB020406, P50 NS091856, P30 DK089503, P30AG053760, UL1TR002240

NSF: 1734853, 1636840, 1416953, 0716055, 1023115

The Elsie Andresen Fiske Research Fund

<http://SOCR.umich.edu>

Collaborators

- **SOCR**: Milen Velez, Alexandr Kalinin, Selvam Palanimalai, Syed Husain, Juana Sanchez, Dennis Pearl, Kyle Siegrist, Rob Gould, Jingshu Xu, Nellie Ponarul, Ming Tang, Asiyah Lin, Nicolas Christou, Hanbo Sun, Tuo Wang, Simeone Marino, Nina Zhou, Yi Zhao, Lu Wang, Qiucheng Wu
- **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 Maciardi, 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

