

SOCR DataSifter: A Statistical Obfuscation Technique enabling Secure & Effective Data Sharing

Ivo D. Dinov

joint work with Nina Zhou, Simeone Marino, Yi Zhao, Lu Wei, Lu Wang

Statistics Online Computational Resource

Health Behavior & Biological Sciences
Computational Medicine & Bioinformatics
Michigan Institute for Data Science
Neuroscience Graduate Program

University of Michigan

<http://SOCR.umich.edu>

Slides Online:
"SOCR News"



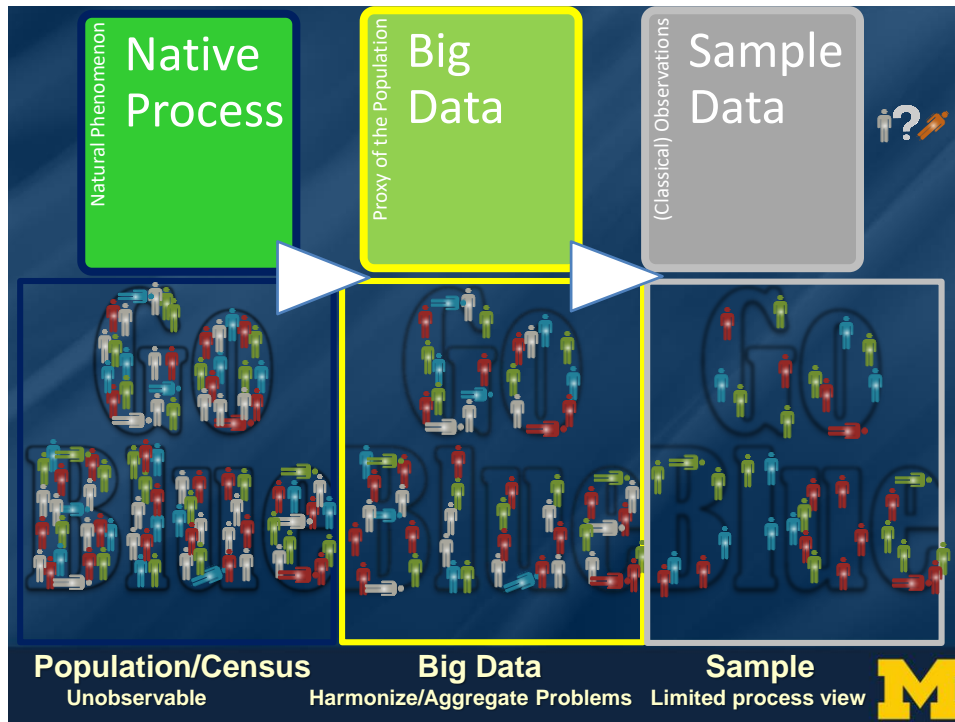
SCHOOL OF NURSING

STATISTICS ONLINE COMPUTATIONAL RESOURCE (SOCR)
UNIVERSITY OF MICHIGAN

Outline

- ☐ Driving biomedical & health challenges
- ☐ Common characteristics of Big Biomedical Data
- ☐ ϵ -Differential Privacy & Fully Homomorphic Encryption
- ☐ *DataSifter: Statistical obfuscation*
- ☐ Case-studies
 - ☐ Applications to Neurodegenerative Disease (PD/AD)
 - ☐ 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

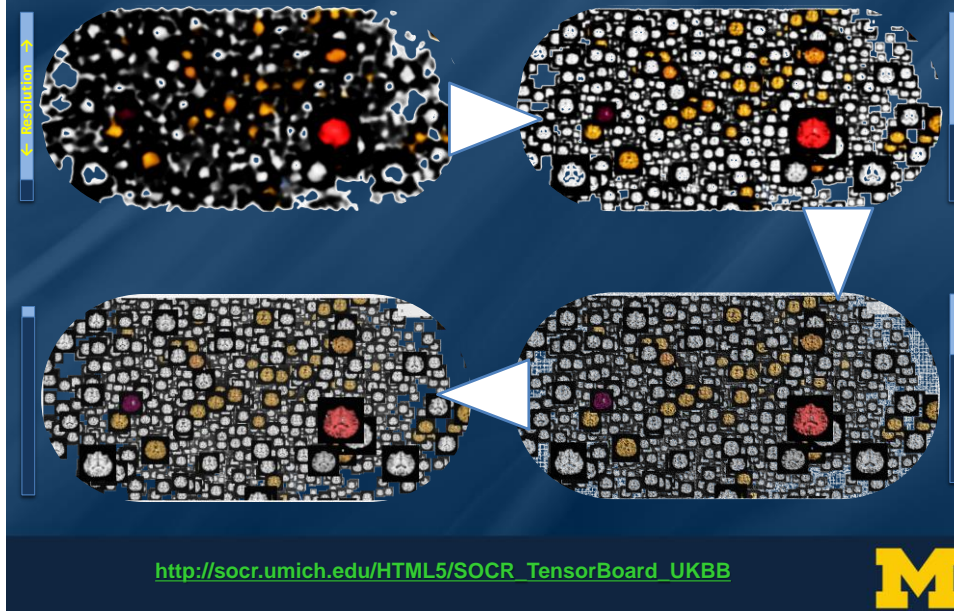
Software developments, student training, service platforms and methodological advances associated with the Big Data Discovery Science all present existing opportunities for learners, educators, researchers, practitioners and policy makers

Dinov (2016) GigaScience

Dinov (2018) Springer

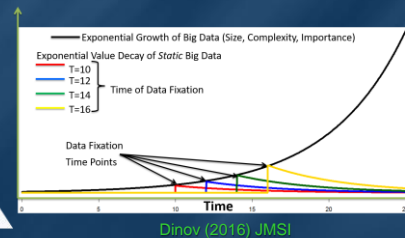
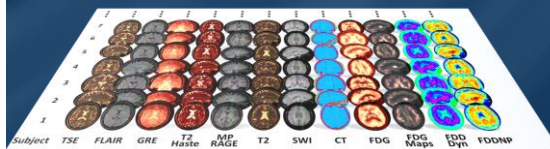


Multiscale/Multimodal NI Data

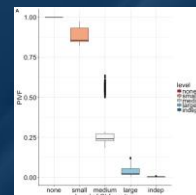


Data Size, Privacy, Usage & Impact

Volume vs. Value of Data



Security vs. Utility



Zhou et al. (2020), pending

ϵ -Differential Privacy (ϵ DP) vs. fully Homomorphic Encryption (fHE)

Category	ϵ DP	fHE
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 derived info
Pros	Theoretical limits on the balance between utility and risk of sharing data	Fast, elegant, and powerful math framework for bijective (encode/decode) encryption
Cons	Difficult for unstructured, skewed, and categorical data	There are limitations on deriving



ϵ -Differential privacy (ϵ DP)

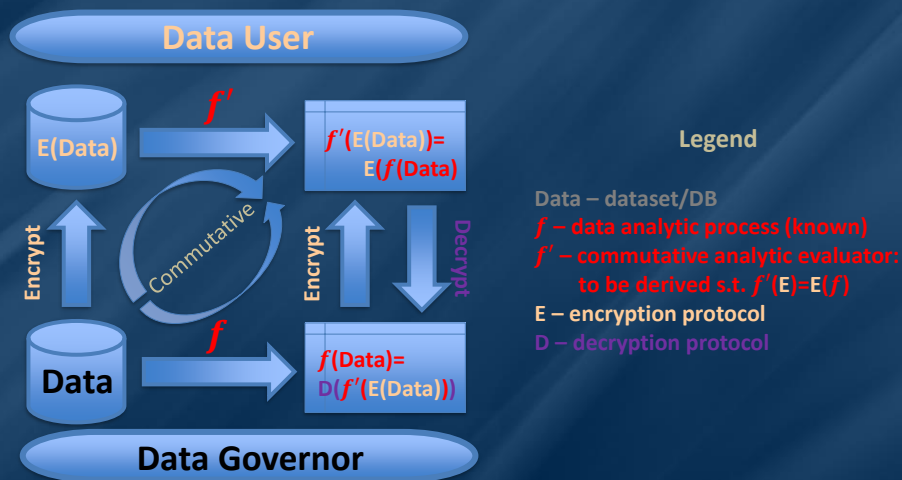
- **Data-features:** $\{C_1, C_2, \dots, C_k\}$, categorical or numerical.
- **DB** = list of cases $\{x_1, x_2, \dots, x_n\}$, $x_i \in \underbrace{C_1 \times C_2 \times \dots \times C_k}_{\text{features}}$, $1 \leq i \leq n$.
- ϵ -Differential privacy relies on adding noise to data to protect the identities of individual records. Given $\epsilon > 0$, **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 (inference), y , the probabilities of correctly guessing D_1 or D_2 knowing y are not significantly different:

$$\frac{P(f(D_1) = y)}{P(f(D_2) = y)} \leq e^\epsilon, \quad \forall y \in \text{Range}(f).$$
- 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 $y = f(D) + (y_1, y_2, \dots, y_m)$, 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

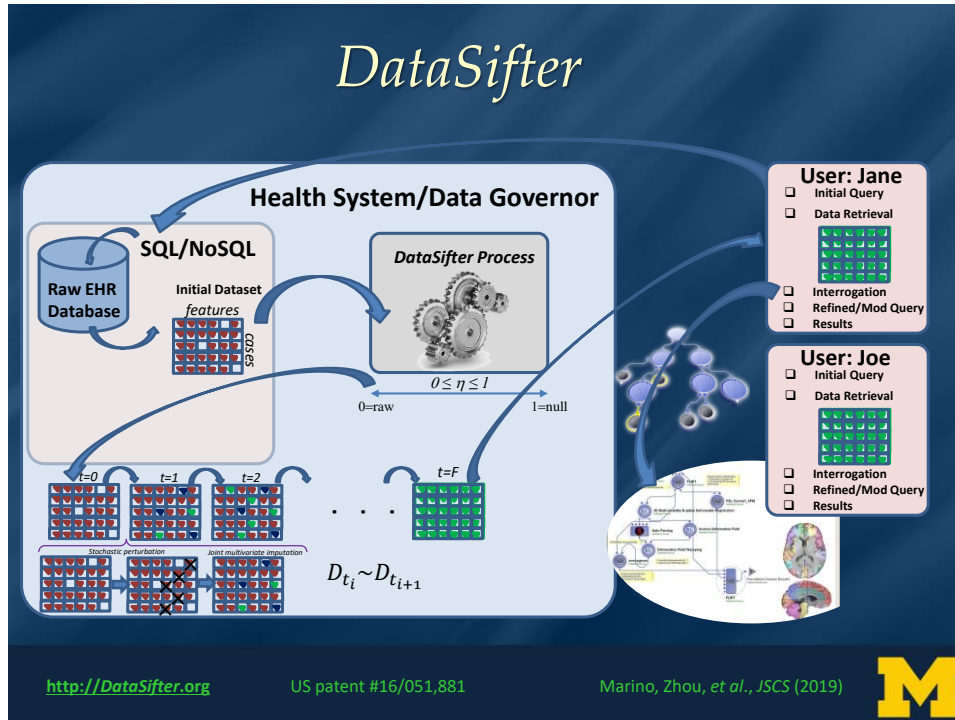
- ❑ DataSifter is an iterative statistical computing approach that provides the data-governors controlled manipulation of the trade-off between sensitive information obfuscation and preservation of the joint distribution.
- ❑ The DataSifter is designed to satisfy data requests from pilot study investigators 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., JSCS (2019)





DataSifter

- To statistically obfuscate the data, DataSifter generates synthetic information and imputes (real or sifter-introduced) missing records by either parametric or semi-parametric prediction models.
- Iterative imputation procedure with (parametric LASSO regularized) Generalized Linear Mixed Model (GLMM)

For each selected time-varying variables: $\mathbf{X}^* = (\mathbf{X}_{1,j}^*, \dots, \mathbf{X}_{n,j}^*)$, fit a prediction model:

$$\eta_{i,j} = g(E(Y_{i,j})) = \mathbf{X}_{i,j}^T \boldsymbol{\beta} + \mathbf{Z}_{i,j}^T \boldsymbol{\gamma}_i$$

where $g(\cdot)$ is a known link function, e.g., logit function for binary data, log function for Poisson count data, etc. $\mathbf{Z}_{i,j}$ is the design matrix of the random effects $\boldsymbol{\gamma}_i \sim N(0, D)$, indexed by $i = 1, \dots, n$ for each subjects and $j = 1, \dots, J_i$ for each time point. Estimate $\boldsymbol{\beta}$ and D using the observed data and impute the missing values by random sampling $\hat{\boldsymbol{\gamma}}_i \sim N(0, \hat{D})$ via best linear unbiased imputation prediction (BLUP): $g^{-1}(\mathbf{X}_{i,j}^T \hat{\boldsymbol{\beta}} + \mathbf{Z}_{i,j}^T \hat{\boldsymbol{\gamma}}_i)$ for $Y_{i,j,k}$ where $(i,j) \in \text{mis}_k$.

- Random Effects-Expectation Maximization tree (RE-EM tree)

Combines the tree-based non-param estimation for fixed effects and parametric estimation for random effects via a linear mixed effect model:

$$Y_{i,j} = f(\mathbf{X}_{i,j,1}, \dots, \mathbf{X}_{i,j,m_k}) + \mathbf{Z}_{i,j}^T \boldsymbol{\gamma}_i + \epsilon_{i,j}, \text{ where } (\epsilon_{1,1}, \dots, \epsilon_{i,J_i})^T \sim N(0, R_i), \text{ and } \boldsymbol{\gamma}_i \sim N(0, D).$$

$f(\cdot)$ is a regression tree and R_i is the variance-covariance structure for i^{th} error term. RE-EM uses the CART tree algorithm to estimate $f(\cdot)$. Assuming we estimated or know $\boldsymbol{\gamma}_i^{(r)}$, the new estimate $\boldsymbol{\gamma}_i^{(r)}$ is obtained by optimizing $Y_{i,j} - \mathbf{Z}_{i,j}^T \boldsymbol{\gamma}_i^{(r)}$. Updating the missing longitudinal variables is achieved iteratively until a stopping criteria is met, e.g., $\frac{\|Y_{\text{mis}_k} - Y_{\text{mis}_k}^{(r)}\|_1}{\|Y_{\text{mis}_k}\|_1} < \epsilon, \forall k = 1, \dots, m_l$.

<http://DataSifter.org> Zhou, et al. (2019) in progress **M**

DataSifter Implementation

Input: Mixed dataset with cross-sectional data and longitudinal data (with/without missing values)

Step 1: Split the data into complete set and missing set for every longitudinal variable $Y_{\cdot,k}$, m_l copies of datasets $\{Y_{obs_k,k}, X_{obs_k,k}, Y_{obs_k,-k}\}$ and $\{Y_{mis_k,k}, X_{mis_k,k}, Y_{mis_k,-k}\}$ for $k = 2, \dots, m_l + 1$.

Step 2: Initiate $Y_{obs_k,-k}^{(0)}$, $\hat{Y}_{mis_k,k}^{(0)}$ and $Y_{mis_k,-k}^{(0)}$ by LOCF, NOCB or mean imputation. Fit logistic regressions for missingness and calculate the probability of being observed for the complete cases of each $Y_{\cdot,k}$.

Step 3: At iteration r^{th} , following variable selection, fit a GLMM LASSO model $f(\cdot)^{(r)}$ on $Y_{obs_k,k}$ with weighted $Y_{obs_k,-k}^{(r)}$ and selected variables $X_{obs_k}^{*(r-1)}$ from X_{obs_k} , $Y_{obs_k,k'r < k}^{(r)}$ and $Y_{obs_k,k'r \geq k}^{(r-1)}$ as possible covariates. Here, $k' < k$ variables are updated in the previous iteration while $k' \geq k$ variables are to be updated. Update $\hat{Y}_{mis_k,k}^{(r)}$ using $X_{mis_k}^{*(r-1)}$ from X_{mis_k} and $Y_{mis_k,-k}^{(r-1)}$ as covariates. Also, update $Y_{obs_k,-k}^{(r)}$, $Y_{mis_k,-k}^{(r)}$ with $f(\cdot)^{(r)}$, for all $k' \neq k$. Check convergence using model predictions for the observed data $Y_{obs_k,k}$ with $f(\cdot)^{(r)}$.

Step 4: Repeat **Step 3** until $\frac{\|Y_{obs_k,k} - \hat{Y}_{obs_k,k}\|_1}{\|Y_{obs_k,k}\|_1} < \epsilon$ or $r = \max_it$. Update using imputed values $Y_{\sum_i J_i \times (m_l + 1)}^*$.

Step 5: Introduce random missingness to m_l longitudinal variables. Keep real values of missing cells as $Y_{mis_k,k}^{**}$.

Step 6: Initiate $Y_{obs_k,-k}^{*(0)}$, $\hat{Y}_{mis_k,k}^{*(0)}$ and $Y_{mis_k,-k}^{*(0)}$ by LOC, NOCB or mean imputation.

Step 7: Use RE-EM or LASSO model $f(\cdot)^{*(r)}$ on $Y_{obs_k,k}^*$ with unweighted $Y_{obs_k,-k}^{*(r)}$ and selected variables $X_{obs_k}^{***(r-1)}$ from X_{obs_k} , $Y_{obs_k,k'r \geq k}^{*(r-1)}$ and $Y_{obs_k,k'r < k}^{*(r)}$ as possible covariates. Update $\hat{Y}_{mis_k,k}^{*(r)}$ using $X_{mis_k}^{***(r-1)}$ from X_{mis_k} and $Y_{mis_k,-k}^{*(r-1)}$ as covariates. Update $Y_{obs_k,-k}^{*(r)}$, $Y_{mis_k,-k}^{*(r)}$ with $f(\cdot)^{*(r)}$, for all $k' \neq k$.

Step 8: Repeat **Step 7** until $\frac{\|Y_{mis_k,k}^* - \hat{Y}_{mis_k,k}^*\|_1}{\|Y_{mis_k,k}^*\|_1} < \epsilon$ or $r = \max_it$. Output the final data $Y_{\sum_i J_i \times (m_l + 1)}^{DS}$ and $X_{n \times m_s}$.

<http://DataSifter.org>

Zhou, et al. (2019) in progress



DataSifter

- ❑ A detailed description and `dataSifter()` R method implementation 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	$0 \leq \eta = \eta_{k_0} + k_1 + k_2 + k_3 + k_4 \leq 1$				
	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

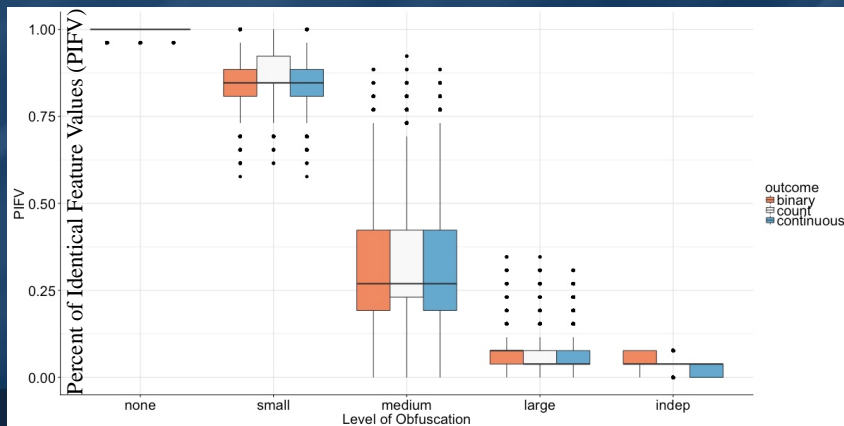
Marino, Zhou, et al., JSCS (2019)



DataSifter Validation

I. Protection of sensitive information (privacy)

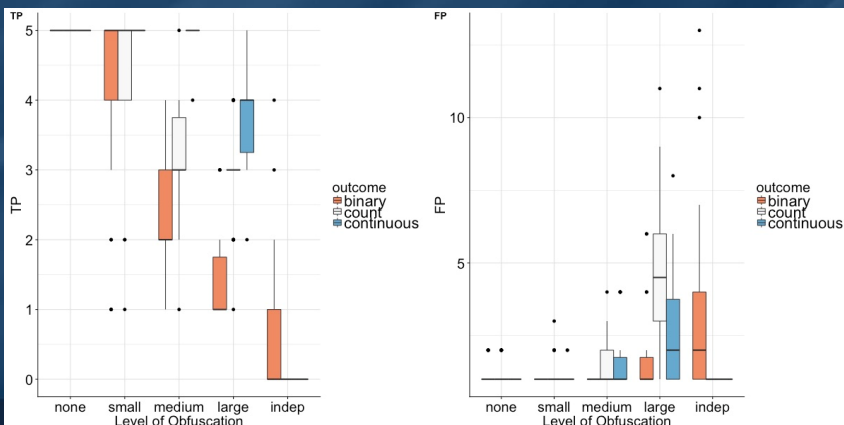
PIFV under Different Privacy Levels. Three simulations are performed using Binary (exp1), Categorical/Count (exp2), and Continuous outcomes (exp3). Each box represents 30 different “sifted” data experiments.



DataSifter Validation

II. Preserving utility information of the original dataset

Logistic Model with Elastic Net Signal Capturing Ability. TP is the number of true *salient features* (total true predictors = 5) captured by the model. FP is the number of *null features* chosen in the model (total null features=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

η	Output	Sex	Age	Acquisition Plane	IQ	thick_std_ctx 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

Autism Brain Imaging Data Exchange (ABIDE) case-study ($n = 1,100$; $k = 2,400$)



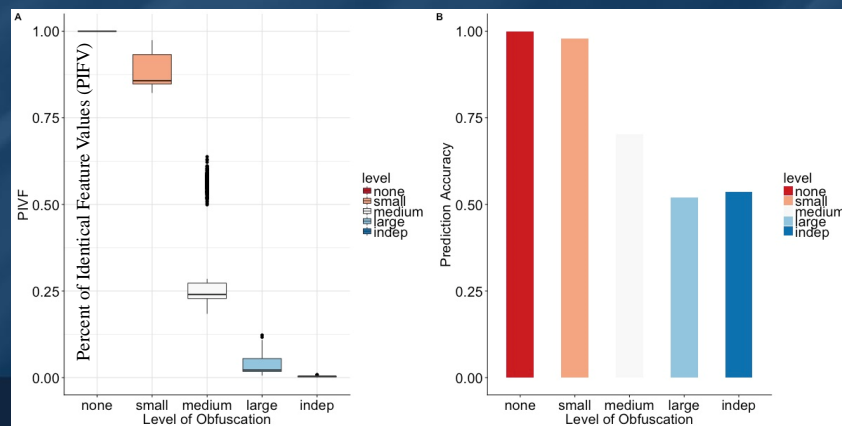
DataSifter Validation

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

PIFVs for ABIDE under different levels of DataSifter obfuscations.

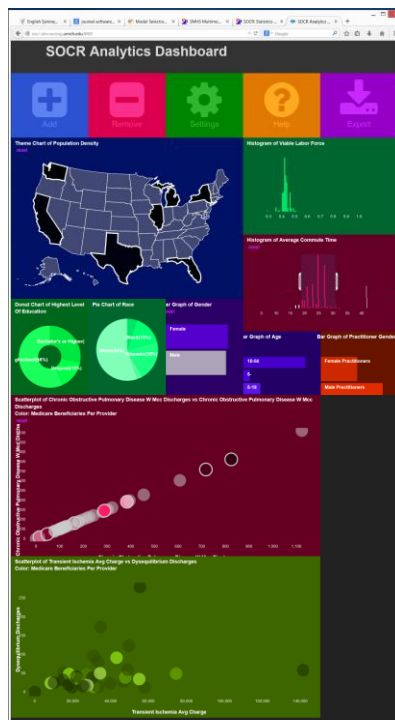
(Left) Each box represents 1,098 subjects among the ABIDE sub-cohort

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



Data Sharing promotes Innovation & Translation

- ❑ SOCR Dashboard
- ❑ Amyotrophic Lateral Sclerosis (ALS, Lou Gehrig's)
- ❑ Neurodegenerative Disorders (Alzheimer's *Parkinson's*)
- ❑ Population epidemiological studies (UKBB)
- ❑ General data integration, augmentation, joining & merging



SOCR Big Data Dashboard

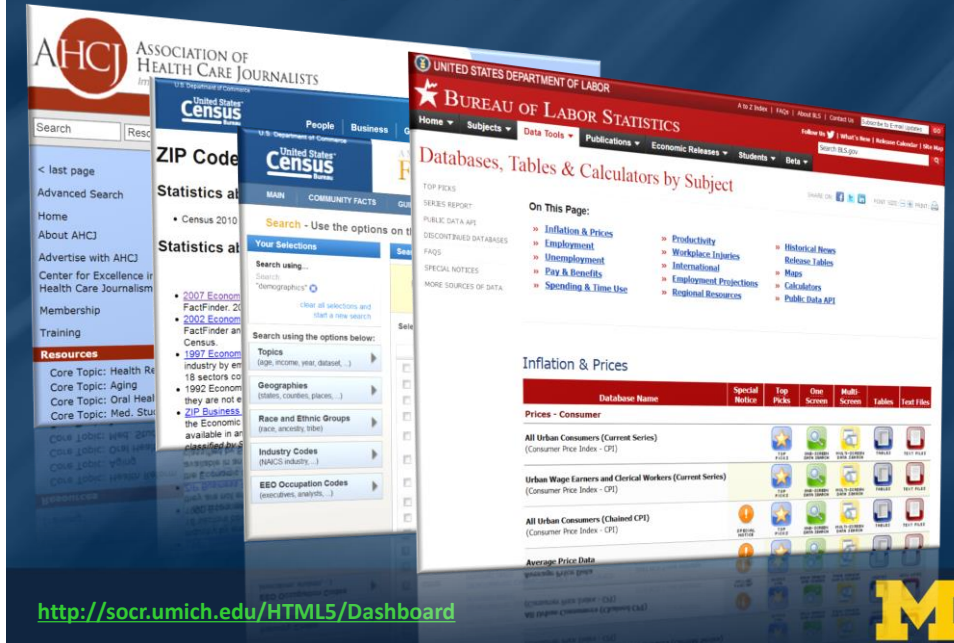
<http://socr.umich.edu/HTML5/Dashboard>

- ❑ 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

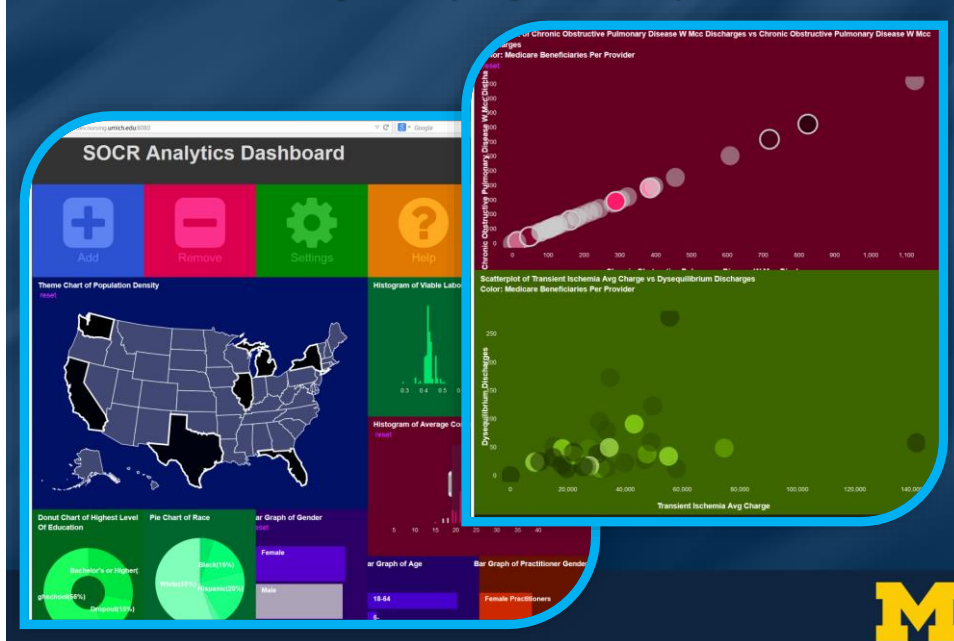
Husain, et al., 2015, PMID:26236573



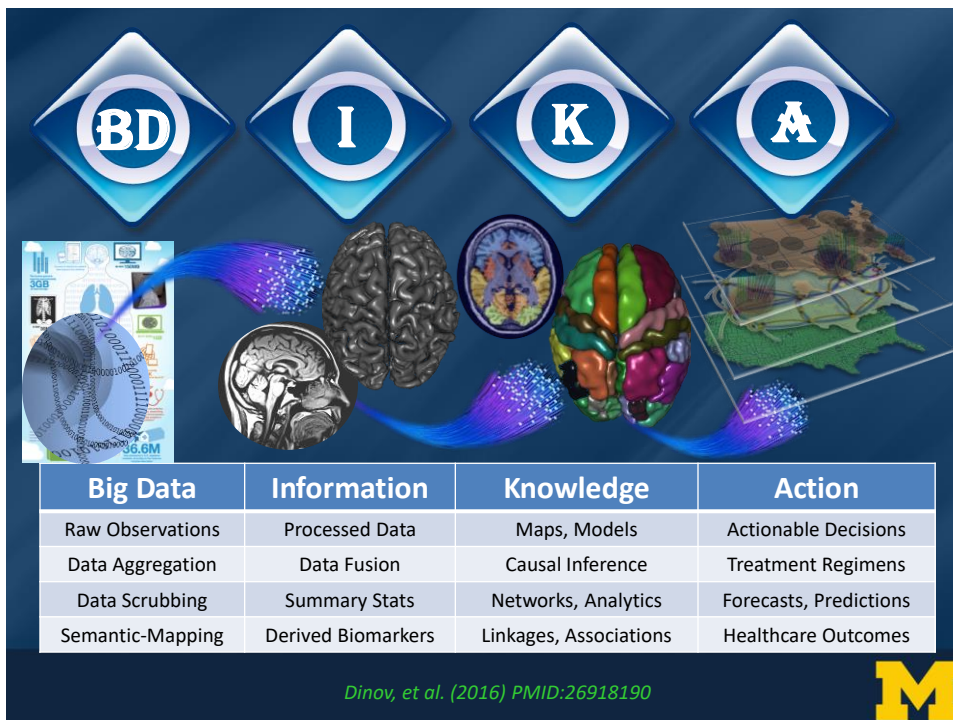
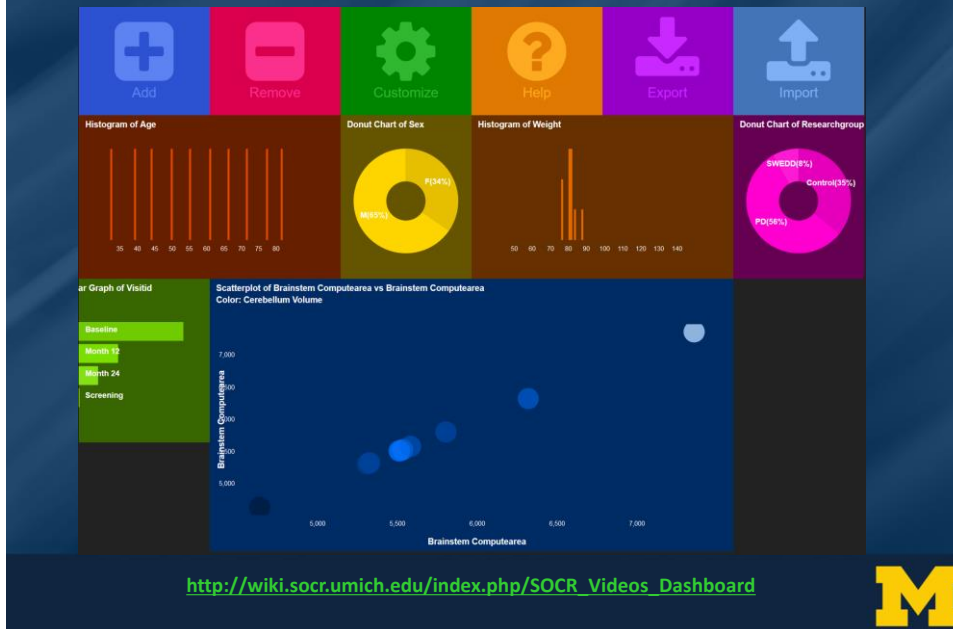
SOCR Dashboard (Exploratory Big Data Analytics): Data Fusion



SOCR Dashboard (Exploratory Big Data Analytics): Associations



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



Why is FAIR Data Sharing Important?

- ❑ Optimum resource utilization (low cost, high efficiency / policy, security, processing complexity)
- ❑ Democratization of the scientific discovery process
- ❑ Enhanced inference (e.g., coverage of rare events, increase of stat power)
- ❑ Increase of Kryder's Law (Data volume) \gg Moore's Law (Compute power)
- ❑ Exponential decay of data-value
- ❑ Incentives innovation, transdisciplinary collaborations, and knowledge dissemination
- ❑ ...

FAIR = Findable + Accessible + Interoperable + Reusable



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: Amyotrophic Lateral Sclerosis 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)

Huang et al. (2017) PLoS

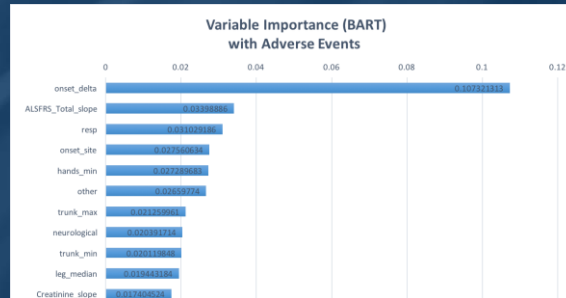
|

Tang, et al. (2018), Neuroinformatics

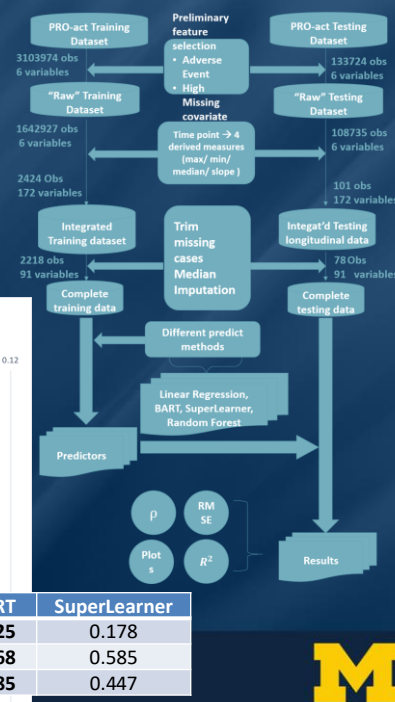


Case-Studies – ALS

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

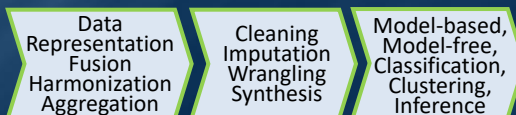


Methods	Linear Regression	Random Forest	BART	SuperLearner
R-squared	0.081	0.174	0.225	0.178
RMSE	0.619	0.587	0.568	0.585
Correlation	0.298	0.434	0.485	0.447



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



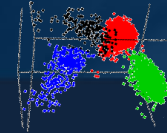
Cluster	Consistency	Variance	Cluster-Size	Silhouette
1	1	0	565	0.58
2	0.986	0.018	427	0.63
3	0.956	0.053	699	0.5
4	0.985	0.018	733	0.5

Tang, et al. (2018), Neuroinformatics

Case-Studies – ALS – Explicating Clustering

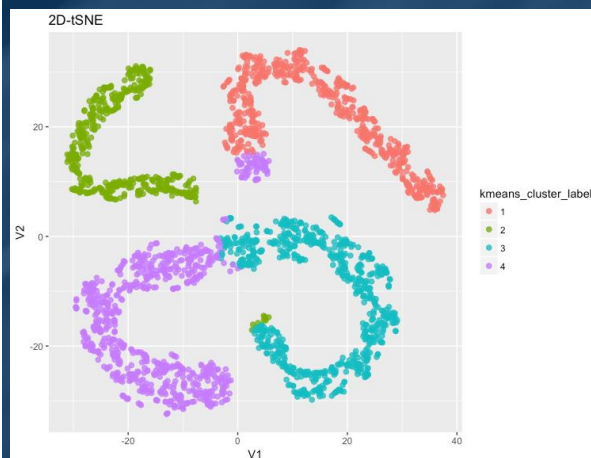
Feature Name	1-2	1-3	1-4	2-3	2-4	3-4
onset_delta.x	1	1	1	1	1	1
onset_delta.y	1	1	1	1	1	1
onset_delta.z	1	1	1	1	1	1
headBloodCath_BSC_min	1	1	1	1	1	1
headBloodCath_BSC_max	1	1	1	1	1	1
Q4_handwriting_max	1	1	1	1	1	1
Q4_handwriting_min	1	1	1	1	1	1
Q4_handwriting_mean	1	1	1	1	1	1
Q9_Climbing_Stairs_max	1	1	1	1	1	1
Q9_Climbing_Stairs_min	1	1	1	1	1	1
Q9_Climbing_Stairs_mean	1	1	1	1	1	1
Q9_Walking_max	1	1	1	1	1	1
Q9_Walking_min	1	1	1	1	1	1
Q9_Walking_mean	1	1	1	1	1	1
touch_max	1	1	1	1	1	1
touch_min	1	1	1	1	1	1
touch_mean	1	1	1	1	1	1
Protein_slope	1	1	1	1	1	1
Creativity_max	1	1	1	1	1	1
Creativity_min	1	1	1	1	1	1
Creativity_mean	1	1	1	1	1	1
creativity_rate_max	1	1	1	1	1	1
head_min	1	1	1	1	1	1
head_max	1	1	1	1	1	1
Q4_Crossing_and_Walking_max	1	1	1	1	1	1
Q4_Crossing_and_Walking_min	1	1	1	1	1	1
Q4_Crossing_and_Walking_mean	1	1	1	1	1	1
Q7_Lumbar_in_Bed_max	1	1	1	1	1	1
Q7_Lumbar_in_Bed_min	1	1	1	1	1	1
Q7_Lumbar_in_Bed_mean	1	1	1	1	1	1
ASPRS_slope	1	1	1	1	1	1
ASPRS_Total_max	1	1	1	1	1	1
ASPRS_Total_min	1	1	1	1	1	1
ASPRS_Total_mean	1	1	1	1	1	1
ASPRS_Total_slope	1	1	1	1	1	1
Neuroticism_max	1	1	1	1	1	1
Neuroticism_min	1	1	1	1	1	1
Neuroticism_mean	1	1	1	1	1	1
leg_min	1	1	1	1	1	1
leg_max	1	1	1	1	1	1
leg_mean	1	1	1	1	1	1
muscle_min	1	1	1	1	1	1
Absolute_Basophil_Count_max	1	1	1	1	1	1
Absolute_Basophil_Count_min	1	1	1	1	1	1
Absolute_Basophil_Count_mean	1	1	1	1	1	1
Absolute_Eosinophil_Count_max	1	1	1	1	1	1
Absolute_Eosinophil_Count_min	1	1	1	1	1	1
Absolute_Eosinophil_Count_mean	1	1	1	1	1	1
Absolute_Lymphocyte_Count_max	1	1	1	1	1	1
Absolute_Lymphocyte_Count_min	1	1	1	1	1	1
Absolute_Lymphocyte_Count_mean	1	1	1	1	1	1
Absolute_Monocyte_Count_max	1	1	1	1	1	1
Absolute_Monocyte_Count_min	1	1	1	1	1	1
Absolute_Monocyte_Count_mean	1	1	1	1	1	1

Feature Name	Between Cluster Significant Differences					
	1-2	1-3	1-4	2-3	2-4	3-4
onset_delta.x	1	1	1	1	1	1
...				...		
Q9_Climbing_Stairs_slope	1			1		
...				...		
leg_max		1	1	1	1	
				...		



Tang, et al. (2018), Neuroinformatics

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), Neuroinformatics



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 forecasts 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 ...

PD phenotypes
Tremor-Dominant (TD)
Postural Instability &
gait difficulty (PI & GD)



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 180 degrees
101	22681	Carotid ultrasound	Maximum carotid IMT (intima-medial thickness) at 210 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 180 degrees
101	22680	Carotid ultrasound	Mean carotid IMT (intima-medial thickness) at 210 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 180 degrees
101	22679	Carotid ultrasound	Minimum carotid IMT (intima-medial thickness) at 210 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 180 degrees

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

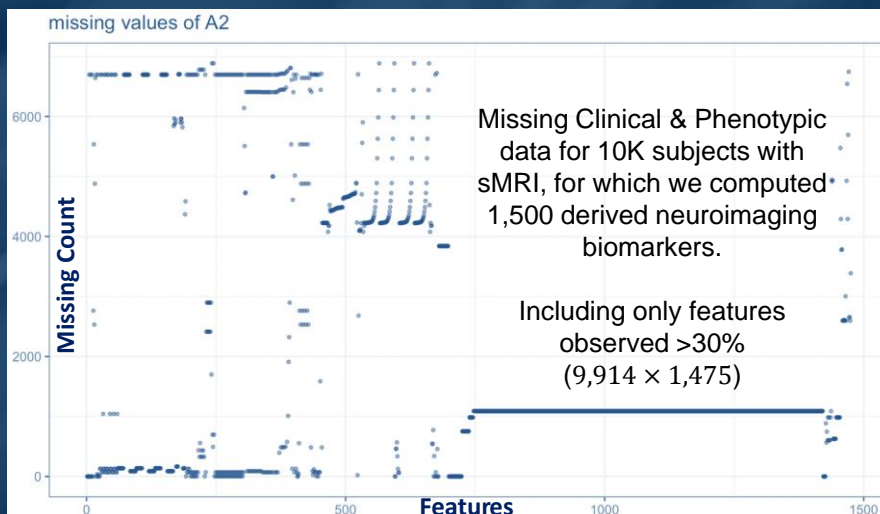
Data Source	Sample Size/Data Type	Summary
UK Biobank	Demographics: > 500K cases Clinical data: > 4K features Imaging data: T1, resting-state fMRI, task fMRI, T2_FLAIR, dMRI, SWI Genetics data	The longitudinal archive of the UK population (NHS)

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

<http://bd2k.org>



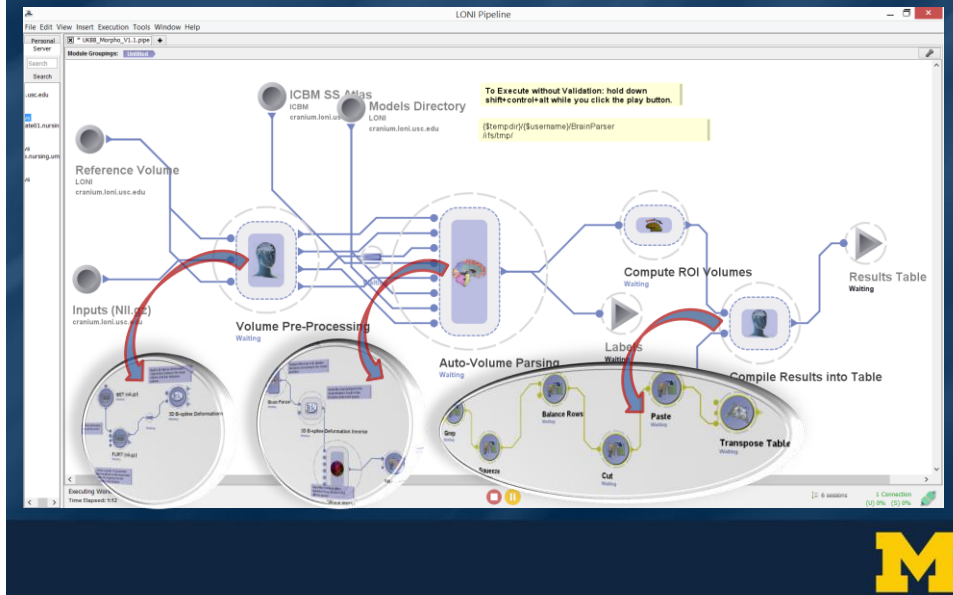
Case-Studies – UK Biobank (Complexities)



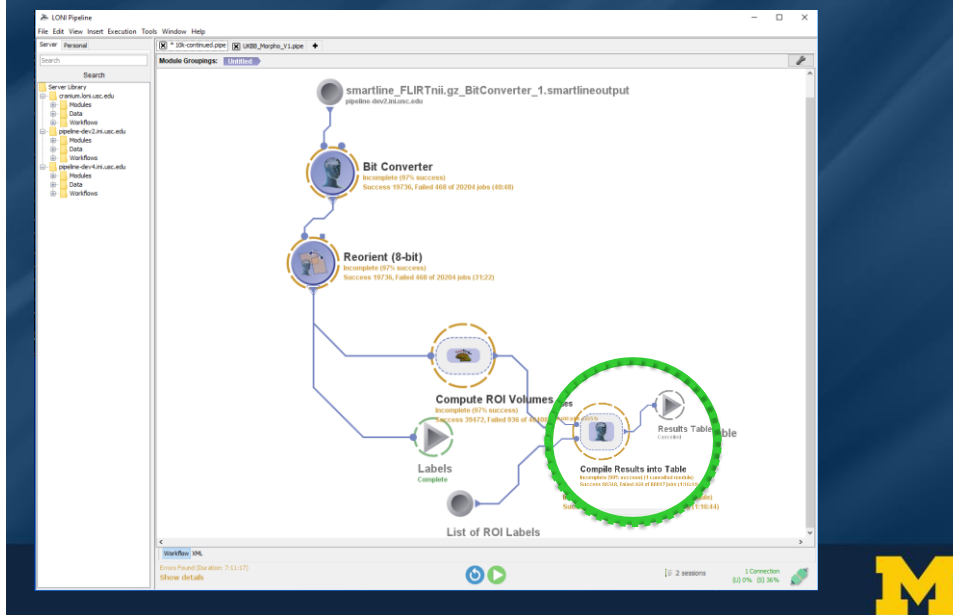
Zhou, et al. (2019), SciRep



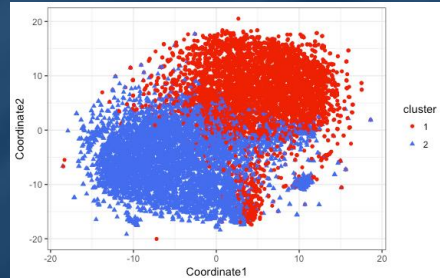
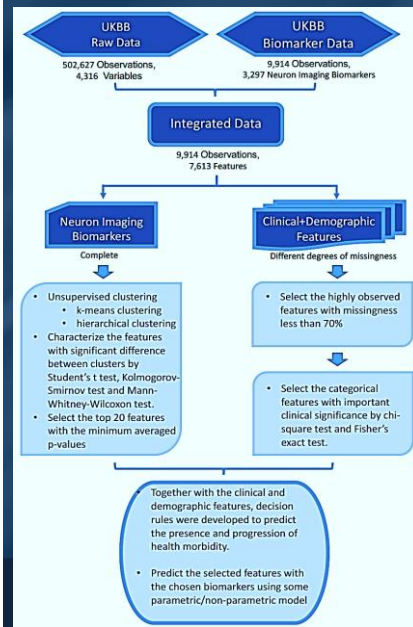
Case-Studies – UK Biobank – NI Biomarkers



Case-Studies – UK Biobank – Successes/Failures



Case-Studies – UK Biobank – Results



k-means clustering

Hierarchical clustering		Cluster 1	Cluster 2
		3768 (38.0%)	528 (5.3%)
	Cluster 1	827 (8.3%)	4791 (48.3%)

Cluster	Consistency	Variance	Cluster-size	Silhouette
1	0.997	0.001	5344	0.09
2	0.934	0.001	4570	0.05



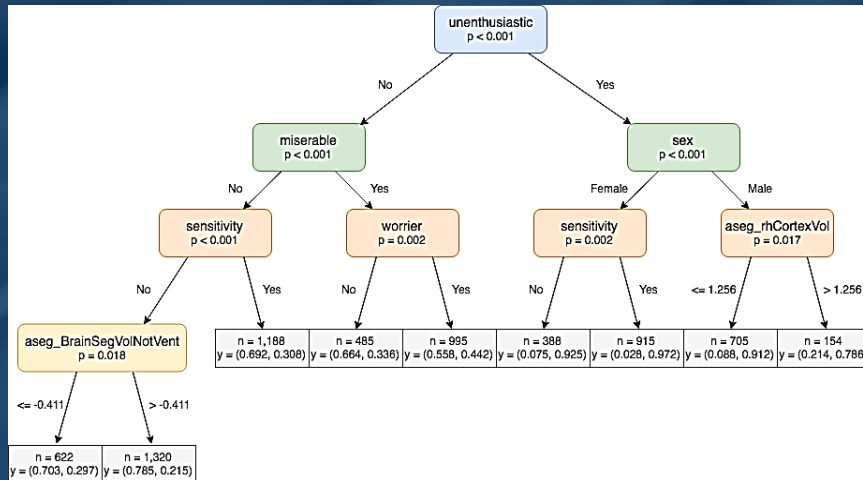
Case-Studies – UK Biobank – Results

Variable	Cluster 1	Cluster 2
Sex		
Female	1,134 (24.7%)	4,062 (76.4%)
Male	3,461 (75.3%)	1,257 (23.6%)
Sensitivity/hurt feelings		
Yes	2,142 (47.9%)	3,023 (56.1%)
No	2,332 (52.1%)	2,151 (43.9%)
Worried/anxious feelings		
Yes	2,173 (48.2%)	2,995 (57.1%)
No	2,337 (51.8%)	2,208 (42.9%)
Risk taking		
Yes	1,378 (31.0%)	1,154 (22.2%)
No	3,084 (69.0%)	3,933 (77.8%)
Guilty feelings		
Yes	1,100 (24.4%)	1,697 (32.1%)
No	3,417 (75.6%)	3,536 (67.9%)
Seen doctor for nerves, anxiety, tension or depression		
Yes	1,341 (29.3%)	1,985 (37.1%)
No	3,237 (70.7%)	3,310 (62.9%)
Alcohol usually taken with meals		
Yes	1,854 (66.7%)	2,519 (76.1%)
No	924 (33.3%)	771 (23.9%)
Drinking		
Yes	1,796 (41.1%)	1,612 (31.1%)
No	2,577 (58.9%)	3,306 (68.9%)
Worry too long after embarrassment		
Yes	1,978 (44.3%)	2,675 (52.1%)
No	2,459 (55.7%)	2,463 (47.9%)
Miserableness		
Yes	1,715 (37.7%)	2,365 (45.1%)
No	2,629 (62.3%)	2,882 (54.9%)
Ever highly irritable/argumentative for 2 days		
Yes	485 (10.7%)	749 (14.5%)
No	4,058 (89.3%)	4,418 (85.5%)
Nervous feelings		
Yes	751 (16.6%)	1,071 (20.8%)
No	3,763 (83.4%)	4,076 (79.2%)
Frequency of tiredness/lethargy in last 2 weeks		
Not at all	2,402 (53.0%)	2,489 (47.8%)
Several days	1,770 (39.0%)	2,127 (40.9%)
More than half the days	187 (4.1%)	300 (5.8%)
Nearly everyday	177 (3.9%)	287 (5.5%)
Alcohol drinker status		
Never	81 (1.8%)	179 (3.4%)
Previous	83 (1.8%)	146 (2.7%)
Current	4,429 (96.4%)	4,992 (93.9%)

Variable	Cluster 1	Cluster 2
Sex		
Female	1,134 (24.7%)	4,062 (76.4%)
Male	3,461 (75.3%)	1,257 (23.6%)
...
Nervous feelings		
Yes	751 (16.6%)	1,071 (20.8%)
No	3,763 (83.4%)	4,076 (79.2%)
...
Frequency of tiredness/lethargy in last 2 weeks		
Not at all	2,402 (53.0%)	2,489 (47.8%)
Several days	1,770 (39.0%)	2,127 (40.9%)
More than half the days	187 (4.1%)	300 (5.8%)
Nearly everyday	177 (3.9%)	287 (5.5%)
Alcohol drinker status		
Never	81 (1.8%)	179 (3.4%)
Previous	83 (1.8%)	146 (2.7%)
Current	4,429 (96.4%)	4,992 (93.9%)



Case-Studies – UK Biobank – Results



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



Acknowledgments

Slides Online:
"SOCR News"

Funding

NIH: P20 NR015331, U54 EB020406, P50 NS091856, P30 DK089503, UL1TR002240, R01CA233487
NSF: 1916425, 1734853, 1636840, 1416953, 0716055, 1023115

Collaborators

- ❑ **SOCR:** Milen Velez, Yongkai Qiu, Zhe Yin, Yufei Yang, Alexandr Kalinin, Selvam Palanimalai, Syed Husain, Matt Leventhal, Ashwini Khare, Rami Elkest, Abhishek Chowdhury, Patrick Tan, Pratyush Pati, Brian Zhang, Juana Sanchez, Dennis Pearl, Kyle Siegrist, Rob Gould, Nicolas Christou, Hanbo Sun, Tuo Wang, **Yi Zhao, Nina Zhou, Yi Wang, Lu Wei, Lu Wang, Simeone Marino**
- ❑ **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 Maciardi, Federica Torri
- ❑ **UMich MIDAS/MNORC/AD/PD Centers:** Cathie Spino, Chuck Burant, Kayvan Najarian, Stephen Goutman, Stephen Strobbe, Hiroko Dodge, Hank Paulson, Bill Dauer, Roger Albin, Chris Monk, Issam El Naqa, HV Jagadish, Brian Atthey



<http://SOCR.umich.edu>