Open Computational Neuroscience: Research, Development & Training



	Outline
о м	otivation
🗆 Pi	llars of Open-Science
🗆 Bi	g Neuroscience
D	ata-Sharing via DataSifter Statistical Obfuscation
	ase-studies
	ALS Study; Parkinson's Disease Study
	Deputation Concurs like Neuroscience (LIKDD)
	Population Census-like Neuroscience (UKBB)



Big Data Characteristics & Challenges

IBM Big Data 4V's: Volume, Variety, Velocity & Veracity					
Big Bio Data Dimensions	Specific Challenges				
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, joint multivariate representation & modeling				
Multi-scale	Interpreting macro \rightarrow meso \rightarrow micro \rightarrow nano scale observations				
Time	Techniques accounting for longitudinal effects (e.g., time corr				
Incomplete	Reliable management of missing data, imputation, obfuscation				

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

Gao et al., SciRep (2018)



Motivation

- Pillars of Open-Science
 - Big Neuroscience
- Data-Sharing via DataSifter Statistical Obfuscation
- Case-studies
- ALS Study; Parkinson's Disease St
- Population Census-like Neuroscience (UKB)
- Spacekime Analytics

Motivation

Pillars of Open-Science

Big Neuroscience

Data-Sharing via DataSifter Statistical Obfuscation

- Case-studies
 - ALS Study; Parkinson's Disease Stu
- Population Census-like Neuroscience (UKBB)
- Spacekime Analytics

From 23 ... to ... 2²³

- Data Science: 1798 vs. 2021
- □ In the 18th century, Henry Cavendish used just 23 observations to answer a fundamental question – "What is the Mass of the Earth?" He estimated very accurately the mean density of the Earth/H₂O (5.483±0.1904 g/cm³)
- □ In the 21st century to achieve the same scientific impact, matching the reliability and the precision of the Cavendish's 18th century prediction, requires a monumental community effort using massive and complex information perhaps on the order of 2²³ bytes
- □ Data Science is about Scalability and Compression 23 → 10M



Motivation Pillars of Open-Science Big Neuroscience Data-Sharing via *DataSifter* Statistical Obfuscation Case-studies ALS Study; Parkinson's Disease Study Population Census-like Neuroscience (UKBB) Spacekime Analytics

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)
- $\hfill\square$ Increase of Kryder's Law (Data volume) \gg Moore's Law (Compute power)
- □ Exponential decay of data-value
- □ Incents innovation, transdisciplinary collaborations, and knowledge dissemination
- ۵ ...

AIR = Findable + Accessible + Interoperable + Reusa



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



ε -Differential privacy (ε DP)

- **Data-features**: { $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_k$ }, categorical or numerical. **DB** = list of cases { x_1, x_2, \dots, x_n }, $x_i \in \underbrace{\mathcal{C}_1 \times \mathcal{C}_2 \times \dots \times \mathcal{C}_k}_{i}$, $1 \le i \le n$.
- $\label{eq:second} \begin{array}{|c|c|c|c|c|} \hline & \varepsilon \mbox{-Differential privacy relies on adding noise to data to protect the identities of individual records. Given $\varepsilon > 0$, algorithm f is $\varepsilon \mbox{-differentially} private if for all possible inputs (datasets/DBs) D_1, D_2 that differ on a single record, and all possible f outputs (inference), y, the probabilities of correctly guessing D_1 or D_2 knowing y are not significantly different: <math display="block">\frac{P(f(D_1)=y)}{P(f(D_2)=y)} \leq e^{\varepsilon}, \qquad \forall y \in Range(f). \end{array}$
- □ 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,
- There are many differentially private algorithms, e.g., random forests, decision trees, k-means clustering, etc.
 E.g., f: D = DB → ℝ^m, the algorithm outputting y = f(D) + (y₁, y₂, ..., y_m)',
- E.g., f: D = DB → ℝ^m, the algorithm outputting y = f(D) + (y₁, y₂, ..., y_m)' with y_i ∈ Laplace (μ = 0, σ = √2 ^{S(p)}/_ℓ), ∀i is ε-differentially private.

Dwork, LNCS, 2008

Homomorphic Encryption (HE)



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 statisticalobfuscation process relies heavily on nonparametric multivariate imputation to preserve the information content of the complex data.

taSifter.org | US patent #10,776,516 | Mari



DataSifter A detailed description and dataSifter() R method implementation are available on our GitHub repository 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 $0 \le \eta = \eta(k_0 + k_1 + k_2 + k_3 + k_4) \le 1$ ion of artificial *k*₄ 0 k₂ 0 None ka: The Small 0.01 0.6 0.8 Medium 0.05 0.2 Large Output synthetic data with independent features Indep





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_ct x .lh.cuneus	curv_ind_ctx _lh_G_front_ inf.Triangul	gaus_curv_ ctx.lh. medialorbitofront al	curv_ind_ctx _lh_S_interm _prim.Jensen
original	Autism	М	31.7	Sagittal	131	0.475	2.1	0.315	NA
none	Autism	м	31.7	Sagittal	131	0.475	2.1	0.315	0.51
small	Autism	М	31.7	Sagittal	131	0.475	2.1	0.315	0.4589
medium	Autism	М	31.7	Sagittal	111	0.548	2.85	0.315	0.463
large	Control	М	18.2	Sagittal	104	0.5347	3.198	0.1625	0.4524
indep	Control	м	15.4	Coronal	104	0.4842	3.383	0.1079	1.002
Au	utism Bra	in Im	aging	Data Exchan	ge (AB	BIDE) case-stu	dy (n = 1,100	; $k = 2,400$)	

DataSifter Validation IV. Clinical Data Application: Using DataSifter to Obfuscate the ABIDE Data IFVS for ABIDE under different levels of DataSifter obfuscations. [Left] Each box represents 1,098 subjects among the ABIDE Sub-cohord [Right] Random forest prediction of binary clinical outcome - autism spectrum disorder (ASD) status (ASD vs. control)

small medium

Percent

Data Sharing Promotes Innovation & Translation

- □ SOCR Dashboard
- □ Amyotrophic Lateral Sclerosis (ALS, Lou Gehrig's disease)
- D Neurodegenerative Disorders (Alzheimer's & Parkinson's)
- Deputation epidemiological studies (UKBB)
- General data integration, augmentation, joining & merging
- Trans-disciplinary education, training, partnerships

Motivation
 Pillars of Open-Science
 Big Neuroscience
 Data-Sharing via DataSifter Statistical Obfuscation
 Case-studies

 ALS Study; Parkinson's Disease Study
 Population Census-like Neuroscience (UKBB)
 Spacekime Analytics

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

ProAct Archive	The time points for all longitudinally varying data elements are aggregated into signatur vectors. This facilitates the modeling and prediction of ALSFRS slope changes over the first three months d (baseline to month 3)	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 Functional</u> Rating Scale (ALSFRS), adverse events, onset_delta, onset_site, drugs use (riluzole). The PRO-ACT training dataset contains clinical and lab test information of <u>8</u> ,635 patients. Information of <u>2</u> ,424 study subjects with valid gold standard ALSFRS slopes used for processing, modeling and analysis
-------------------	--	---



Case-Studies – ALS

- <u>Main Finding</u>: 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

uet
Silha
0.58
0.63
0.5
M



Case-Studies – ALS – **Dimensionality Reduction** 2D t-SNE Manifold embedding



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 <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).
- <u>Results</u>: 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%

Cas	e-St	udie	es – .	Park	insoi	n's Disease
PD_Subtype 60- 40- 20-	Tremor_score	PIGD_score	patSpeed_Off	MoCA	Pg. Surge	
	Ĺ	Cor : 0.065 0: 0.0305 1: 0.112	Cor : -0.17 0: -0.146 1: -0.242	Cor : -0.111 0: -0.0756 1: -0.197	Turne_son	
		h	Cor : -0.644 00.493 1: -0.668	Cor : -0.108 (-0.00449 1: -0.111	PICD_score	Falls in PD are extremely difficult to predict
15 ⁻ 10 ⁻ 1 04 ⁻ 35 ⁻				Cor : 0.215 0: 0.207 1: 0.113	set Speed, Off	PD phenotypes
20-						Postural Instability & gait difficulty (PI & GD)
20- 0- 40- 20- 0 N PI TD				الالله. مارالي . في في في في	a i	M

Case-Studies – Parkinson's Disease							
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>	0.683	0.871	0.778	0.806	2.677	<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, partIL_sum, Attention, DGI, FOG_O, 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

- \circ 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 – NI Biomarkers







Case	-Stu	die	s -	- UK Biobank –	Result	S
Yarabie Sex	Chater 1	Course				
Peeralle Male	1,134(26.7%) 3,461(75.3%)	4,062(NL4N) 1,257(23.6N)				
Sensitivity/hortfeelings Tes	2,142(47.9%)	3,023 (NR-6N)				
No Wonter/andous feelings	2,832(52.1N)	2,151(£1.6N)				
Tes No	2,178(d8.2%) 2,887(51.8%)	2,995(\$7.4%) 2,308(\$2.4%)		Variable	Cluster 1	Cluster 2
Nex Laborg	1,878(81.0N)	1,154(22.7%)		variable	Cluster 1	Cluster 2
Guilty feelings		1,011(77.3%)		Sex		
	3,417(75.6N)	3,536(67.6N)		Female	1,134 (24.7%)	4,062 (76.4%)
Tes	1,841(29.8N)	1,885(87.5N)		Male	3,461 (75,3%)	1,257 (23.6%)
Alizabel seaaby takes with resals	1,854(66.7%)	7.510(12.510)			-,(,	-,
No.	936 [88.36]	771(23.4%)				
	1,796(41.1N) 3.577058.000	1,652(83.3%) 1,876(86.7%)		Normania facilizar		
Wony too long after embasticoment TH	1.978(44.3%)	2.675(02.13)		Ner Yous reenings	754 (46 600)	4.074 (20.00()
No.	2,491(55.7%)	2,462(67.9%)		Yes	/51 (16.6%)	1,071 (20.8%)
	1,715(87.7%)	2,365(05.33)		No	3,763 (83.4%)	4,076 (79.2%)
Free highly initialite/argumentative for 2 days	495 (10.76)	100114 510		L.,		
No Netwice feelings	4,038(89.8N)	4,418(85.5%)		···	•••	
	751(16.6%) 3.753(83.4%)	1,071(20.8%) 4,076(29.2%)		Frequency of tiredness/lethargy in		
Ever depressed for a whole week Tes	2.126(48.15)	2.709(52.9%)		last 2 weeks	2.402 (53.0%)	2.489 (47.8%)
No. Note: The second for a whole week	2,847(51.9%)	2,438(47.1%)		Not at all	1 770 (20.0%)	2,127 (40.9%)
	1,346(80.3%)	1,768(84.3%)		Notatali	1,770 (59.0%)	2,127 (40.9%)
Steeplers/meannia Menus (Cashe	1.007/20.000	1 181 (72 190		Several days	187 (4.1%1)	300 (5.8%)
	2,302(47.9%)	2,571(68.6N) 1,552(78.6N)		More than half the days	177 (3.9%)	287 (5.5%)
Getting up in morning Not at all easy	189(8.1%)	20916.7%)		Nearly everyday		
Not very skip	538(11.9%)	880(15.8%)		Alcohol drinker status		
Very easy New Germander	1,526(88.7%)	1,505(28.7%)		Never	81 (1.8%)	179 (3.4%)
Never/Sarry SomeTimes	2,487(54.5N) 1,774(88.8N)	8,238(61.5N) 1,798(04.2N)		Brovious	92 (1 9%)	146 (2.7%)
Steady Presence of Strategy Systems in 2 works	307 (6.7%)	228(4.3%)		Connect	4 420 (00 40()	140 (2.7%)
Not at all Several days	2,422(58.0N) 1,770(89.0N)	2,489(67.8%) 2,127(02.9%)		Current	4,429 (96.4%)	4,992 (93.9%)
More than half the days Nearly everyday	187(4.1N1) 177(3.9%)	800 (5.8%) 287 (5.5%)				
Alushid drinker status Névér	81(1.8%)	179(8.4%)				
Previous Customet	83(5.8%) 6.629(96.6%)	106(2.7%)				



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						

¥.



Mathematical-Physics \Longrightarrow Data/Neuro Sciences

Mathematical-Physics

A <u>particle</u> is a small localized object that permits observations and characterization or its physical or chemical properties An <u>observable</u> a dynamic variable about particles that can be measured Particle <u>stude</u> is an observable particle characteristic (e.g., position, momentum) Particle <u>stude</u> is a colsection of independent particles and observable characteristics, in a closed system <u>Wave-function</u>

Reference-Frame transforms (e.g., Lorentz) State of a system is an observed measurement of all particles – wavefunction A particle system is computable if (1) the entire system is logical, consistent, complete and (2) the unknown internal states of the system don't influence the computation (wavefunction, intervals, probabilities, etc.)

Data/Neuro Sciences
An object is something that exists by itself, actually o potentially, concretely or abstractly, physically or incorporeal (e.g., person, subject, etc.)
A <u>feature</u> is a dynamic variable or an attribute about a object that can be measured
<u>Datum</u> is an observed quantitative or qualitative value an instantiation, of a feature
Problem, aka Data System, is a collection of independent objects and features, without necessarily being associated with apriori hypotheses
Inference-function
Data transformations (e.g., wrangling, log-transform
Dataset (data) is an observed instance of a set of datum elements about the problem system, $O = \{X, Y\}$
Computable data object is a very special representation of a dataset which allows direct

application of computational processing, modeling, analytics, or inference based on the observed dataset

Dinov & Velev, De Gruyter (2021)

<section-header><section-header>Spacekine Analytics: fMRI Examplea structure Reconstruction of (space=2, time=1) MRI signala structure Reconstruction of (space=2, time=1) MRI signala structure reconstruction of (space=2, time=1) MRI signalb structure reconstruction using trial
b spacestructure reconstruction using
trians-angle; kime=time=(magnitude, place)b structure reconstruction
b spacec structure reconstruction
spacec structure reconstruction
spacec structure reconstruction
spacec structure reconstruction
spacec structure reconstruction
spacec structure reconstruction
spacec structure reconstructionc structure reconstruction
spacec structure reconstructure
spacec structure reconstructurec structure reconstructure
spacec structure reconstructurec structure reconstructurec structure reconstructure
spacec structure reconstructurec structure reconstructure<



In the 5D spacekime manifold, time-series curves extend to kime-series, i.e., surfaces parameterized by kime-magnitude (t) and the kime-phase (φ).

Kime-phase aggregating operators that can be used to transform standard time-series curves to spacekime kimesurfaces, which can be modeled, interpreted, and predicted using advanced spacekime analytics.









Spacekime Analytics: Demos

Tutorials

- https://socr.umich.edu/HTML5/SOCR_TensorBoard_UKBB
 https://DSPA.predictive.space
 https://TCIU.predictive.space | https://SpaceKime.org

🖵 R Package

GitHub

https://github.com/SOCR/TCIU

Summary

- Big Neuroscience Challenges
- □ Open-Science Drivers

Data-Sharing via DataSifter Statistical Obfuscation

Case-studies

- □ ALS Study; Parkinson's Disease Study
- Deputation Census-like Neuroscience (UKBB)
- □ Spacekime Analytics

