# *DataSifter*: Sharing of Sensitive Data 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"*

**M** | STATISTICS ONLINE COMPUTATIONAL RESOURCE (SOCR)
UNIVERSITY OF MICHIGAN

---

# My Perspective on Stats & Math @ FSU

❑ 1947 – FSU
  ❑ FL State College(1905)←FL University(1883)←West FL Seminary (1857)
❑ 1906 – Math (Chair: Elmer Smith; til' 1942, cf. Smith Hall)
❑ 1950's ➔Topology (Morton Curtis, Thomas Wade, Orville Harrold)
❑ 1959 – Stats spins off of Math
❑ 1967 – De Witt Sumners (RO Lawton Prof.) joins
❑ 1982 – Fred Huffer joins Stats (from Stanford)
❑ 1993-1998 – ID dually enrolled in Math/Stats
❑ …
❑ 2019 – FSU-Stats 60th anniversary

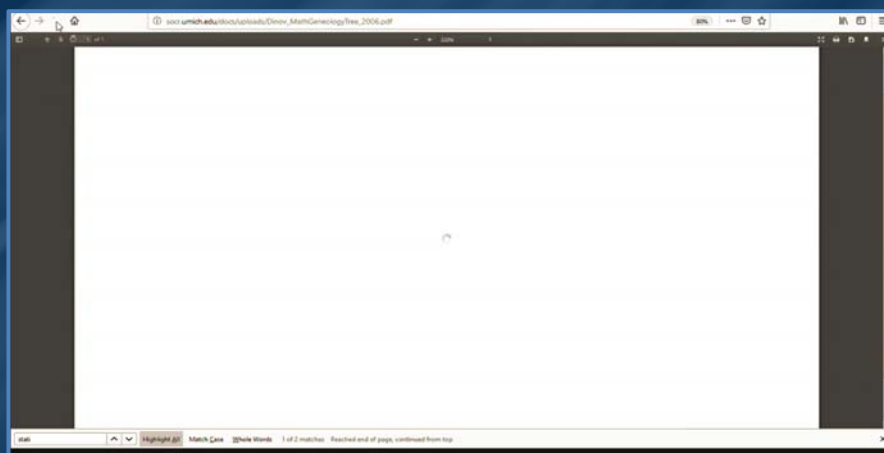**https://www.math.fsu.edu/News/Math_Newsletter_Fal2001.pdf**
**https://www.math.fsu.edu/News/Math_Newsletter_Spr2006.pdf**
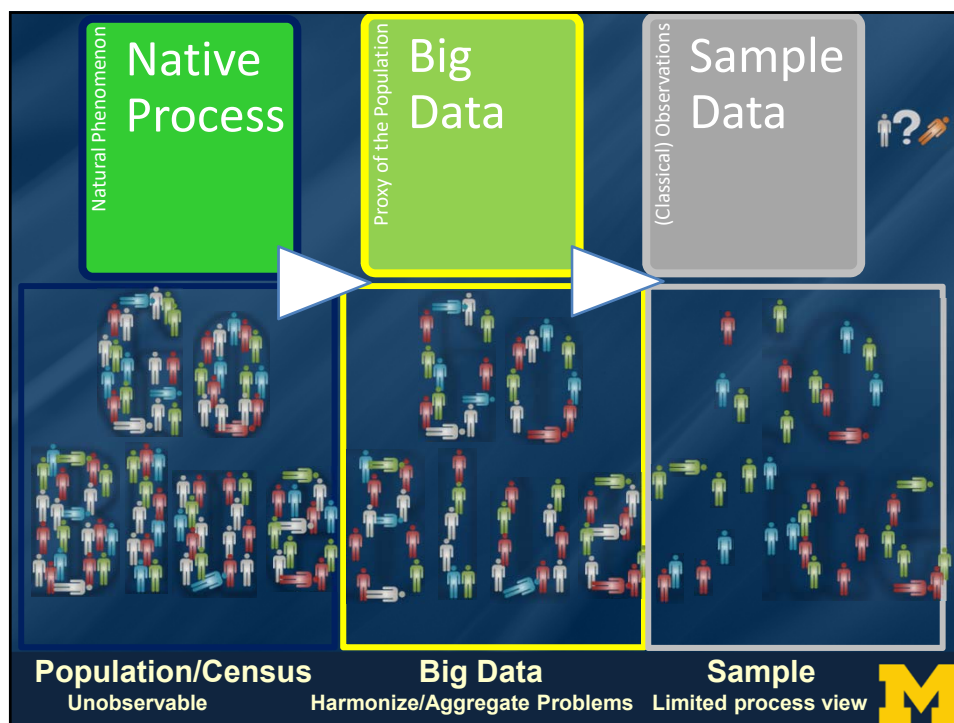
# FSU Stats Alumni (Academic) Pedigree

# Outline

❑ Driving biomedical & health challenges

❑ Common characteristics of Big Neuroscience Data

❑ $\varepsilon$-Differential Privacy & 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 (2016) GigaScience          Dinov (2018) Springer

# $\varepsilon$-Differential Privacy ($\varepsilon$DP) vs. fully Homomorphic Encryption (fHE)

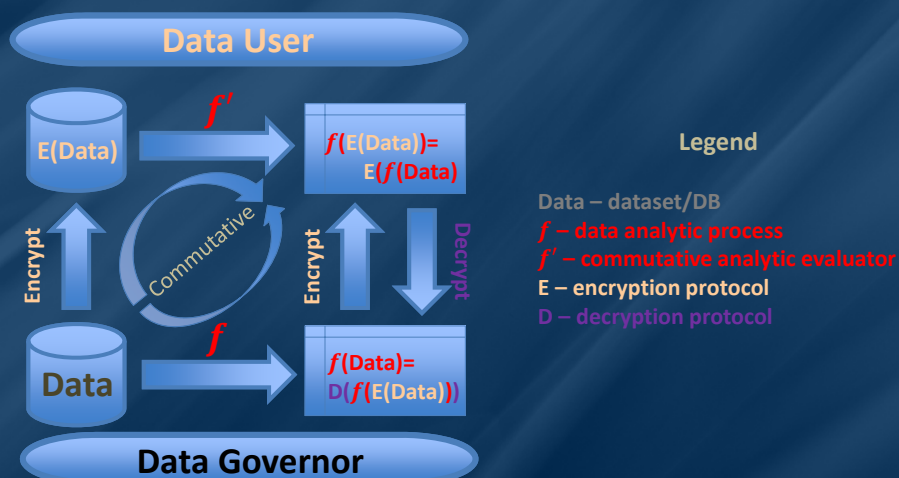| Category | $\varepsilon$DP | fHE |
|---|---|---|
| Goal | Mine information in a DB without compromising privacy; no access to inspect individual DB entries | Provide a secure encryption allowing program execution on encrypted data; encrypt results, interpretation requires ability to decrypt the data |
| Pros | Theoretical limits on the balance between utility and risk of sharing data | Elegant and powerful math framework for bijective (encode/decode) encryption. Fast |
| Cons | Difficult for unstructured, skewed, and categorical data | There are limitations on deriving $f'$ – commutative analytic evaluators |

# $\varepsilon$-Differential privacy ($\varepsilon$DP)

❑ **Data-features**: $\{C_1, C_2, \ldots, C_k\}$, categorical or numerical.
❑ **DB** = list of cases $\{x_1, x_2, \ldots, x_n\}$, $x_i \in C_1 \times C_2 \times \cdots \times C_k$, $1 \le i \le n$.

❑ $\varepsilon$-Differential privacy relies on adding noise to data to protect the identities of individual records. An **algorithm $f$** is $\varepsilon$-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 $D_2$:

$$\frac{P(f(D_1) = y)}{P(f(D_2) = y)} \le e^{\varepsilon}, \qquad \forall y \in Range(f).$$

❑ 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 $\|f(D_1) - f(D_2)\|_1 \le S(f)$
❑ There are many differentially private algorithms, e.g., random forests, decision trees, k-means clustering, etc.
❑ E.g., $f: D = DB \to R^m$, the algorithm outputting $f(D) + (y_1, y_2, \ldots, y_m)$, with $y_i \in Laplace\left(\mu = 0, \sigma = \sqrt{2}\frac{S(f)}{\varepsilon}\right), \forall i$ is $\varepsilon$-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, *et al.*, *JSCS* (2019)

---

# 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.

$k_0$: 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

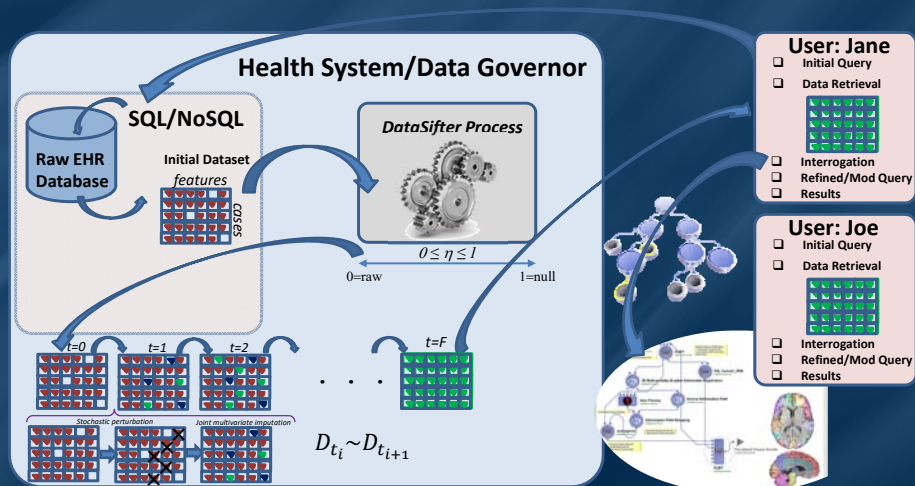$k_4$: The fraction of closest subjects to be considered as neighbours of a given subject

| Obfuscation level | $0 \le \eta = \eta(k_0 + k_1 + k_2 + k_3 + k_4) \le 1$ | | | | |
|---|---|---|---|---|---|
| | $k_0$ | $k_1$ | $k_2$ | $k_3$ | $k_4$ |
| 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 | | | | |

http://*DataSifter*.org   US patent #16/051,881   Marino, *et al., JSCS* (2019)

# DataSifter



**Health System/Data Governor**

**SQL/NoSQL**

**Raw EHR Database**   **Initial Dataset** *features*   *cases*

**DataSifter Process**

$0 \le \eta \le 1$

0=raw   1=null

$t=0$   $t=1$   $t=2$   $t=F$

*Stochastic perturbation*   *Joint multivariate imputation*

$D_{t_i} \sim D_{t_{i+1}}$

**User: Jane**
- Initial Query
- Data Retrieval
- Interrogation
- Refined/Mod Query
- Results

**User: Joe**
- Initial Query
- Data Retrieval
- Interrogation
- Refined/Mod Query
- Results

http://*DataSifter*.org   US patent #16/051,881   Marino, *et al., JSCS* (2019)
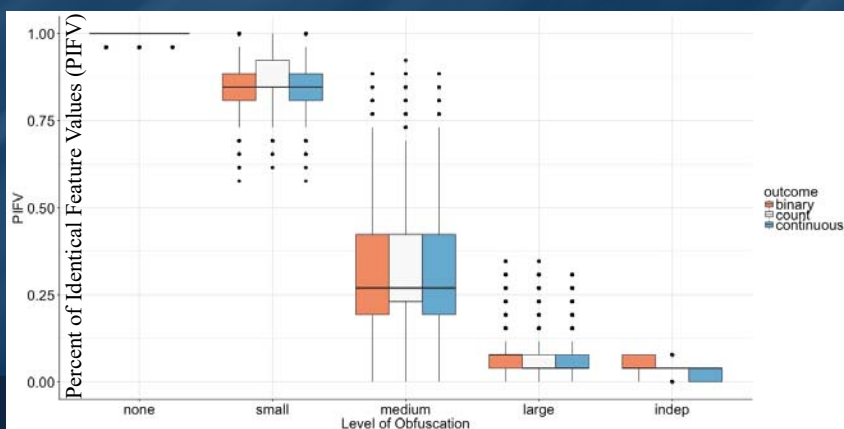
# *DataSifter* Validation

**I. Protection of sensitive information (privacy)**
PIFV under Different Privacy Levels. Binary outcome refers to the first experiment; Count refers to the second experiment; Continuous refers to the third experiment.
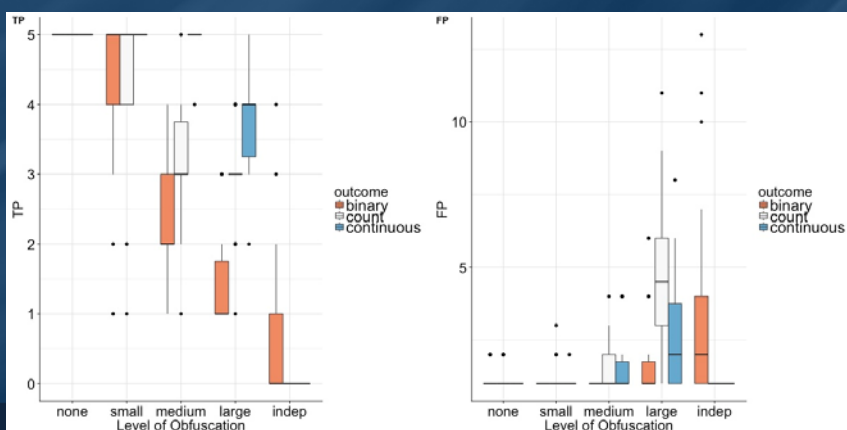Each box represents 30 different "sifted" data or 30,000 "sifted" cases.



# *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).



7

# *DataSifter* Validation

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

Comparing the Original and "Sifted" Data for a random (22nd) ABIDE Subject

| $\eta$ | 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 | 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 |

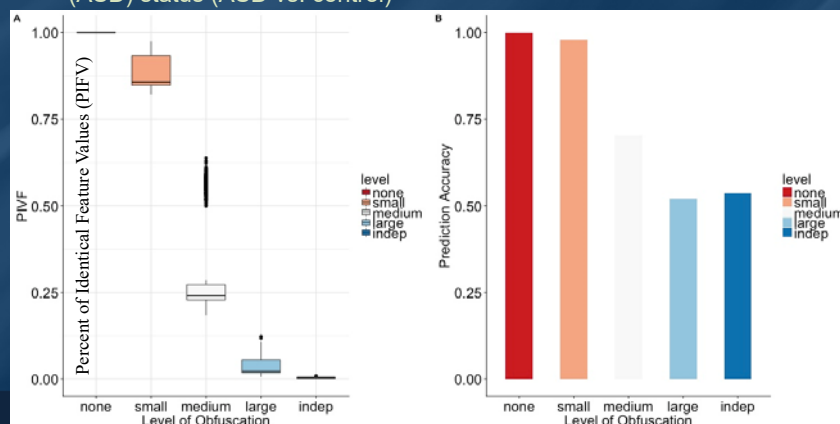Autism Brain Imaging Data Exchange (ABIDE) case-study

# *DataSifter* Validation

**IV. Clinical Data Application: Using DataSifter to Obfuscate ABIDE**
    PIFVs for ABIDE under different levels of DataSifter obfuscations.
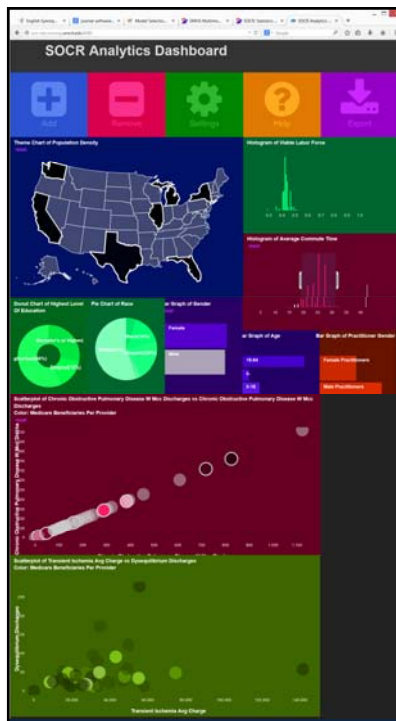    Each box represents 1,098 subjects among the ABIDE sub-cohort
    Random forest prediction of binary clinical outcome - autism spectrum disorder
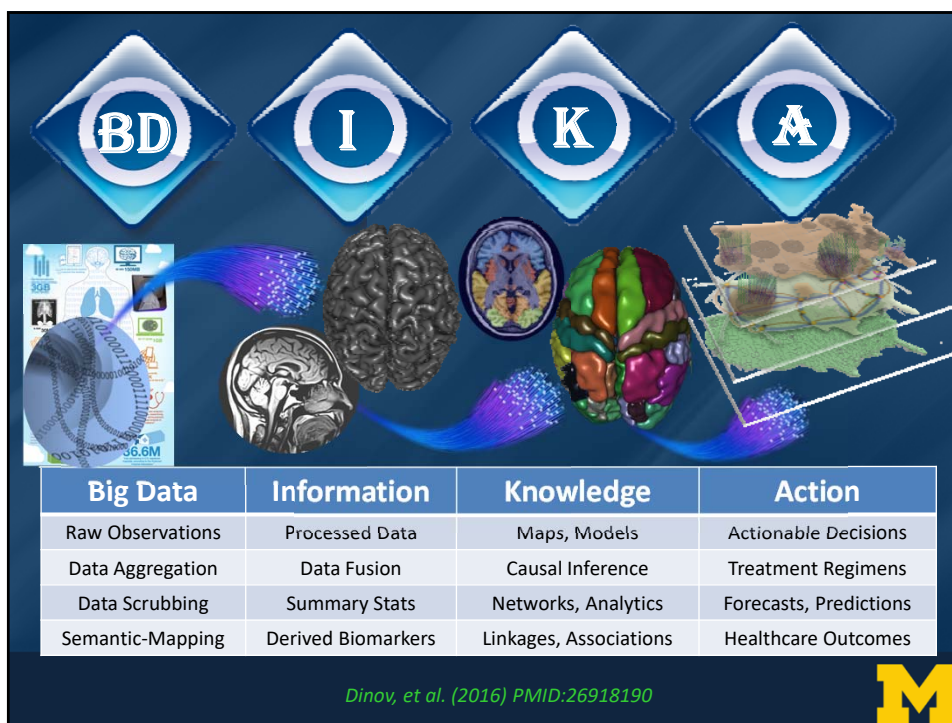    (ASD) status (ASD vs. control)

Raw Data → Statistical Obfuscation → Data Analytics
*Sensitive Info*     *DataSifter*     *Model−based*
                                      *Model−Free*

APPLICATIONS



# 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*

| Big Data | Information | Knowledge | Action |
|---|---|---|---|
| Raw Observations | Processed Data | Maps, Models | Actionable Decisions |
| Data Aggregation | Data Fusion | Causal Inference | Treatment Regimens |
| Data Scrubbing | Summary Stats | Networks, Analytics | Forecasts, Predictions |
| Semantic-Mapping | Derived Biomarkers | Linkages, Associations | Healthcare Outcomes |

*Dinov, et al. (2016) PMID:26918190*

# 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) >> Moore's Law (Compute power)

❑ Exponential decay of data-value

❑ Incents innovation, transdisciplinary collaborations, and knowledge dissemination

❑ …

FAIR = Findable + Accessible + Interoperable + Reusable

# 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*
- **<u>Classification vs. Clustering</u>**:
  - 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

Data Representation Fusion Harmonization Aggregation → Cleaning Imputation Wrangling Synthesis → Model-based, Model-free, Classification, Clustering, Inference

| 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 | 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-tSNE

kmeans_cluster_label
- 1
- 2
- 3
- 4

**2D t-SNE Manifold embedding**

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

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

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

*Tang, et al. (2018), Neuroinformatics*

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

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

---

# Acknowledgments

Slides Online: "SOCR News"

US patent #16/051,881