**Slide 1**

# Data Sharing + Open, Rigorous & Reproducible Science

## Ivo D. Dinov

**Statistics Online Computational Resource**
Health Behavior & Biological Sciences
Computational Medicine & Bioinformatics

University of Michigan

https://SOCR.umich.edu

*Slides Online: "SOCR News"*

STATISTICS ONLINE COMPUTATIONAL RESOURCE (SOCR)
UNIVERSITY OF MICHIGAN

1

**Slide 2**

# Outline

- ❑ Pillars of Open-Science
- ❑ Rationale (Pros & Cons)
- ❑ Big Data Sharing
- ❑ *DataSifter: Statistical obfuscation*
- ❑ Case-studies
  - ❑ ALS Study
  - ❑ Population Census-like Neuroscience (UKBB)
  - ❑ Spacekime Analytics

2

**Slide 3**

# Pillars of Open Data Science (HS650 / Bioinfo501)



3

**Slide 4**

# Sources: 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 (2023) Springer

4

**Slide 5**



Native Process (Natural Phenomenon)

Big Data (Proxy of the Population)

Sample Data ((Classical) Observations)

**Population/Census** Unobservable

**Big Data** Harmonize/Aggregate Problems

**Sample** Limited process view
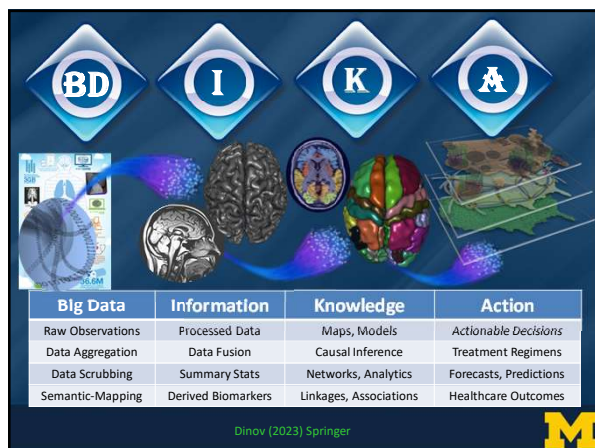
5

**Slide 6**

# From 23 … to … $2^{23}$

- ❑ Data Science: 1798 vs. 2024
- ❑ 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/$H_2O$ (5.483±0.1904 g/cm$^3$)
- ❑ 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 $2^{23}$ bytes
- ❑ Data & Information Science $\cong$ Scalability & Compression (per Gerald Friedland/Berkeley): 23 ➔ $2^{23} \cong$ 10M

Cavendish (1798) Philosophical Transactions of the Royal Society of London    |    Dinov (2016) JSMI

6

## Slide 7



| 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 |

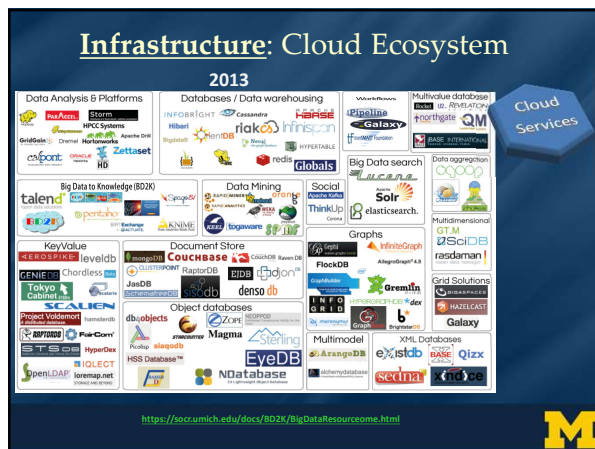Dinov (2023) Springer

7

## Slide 8
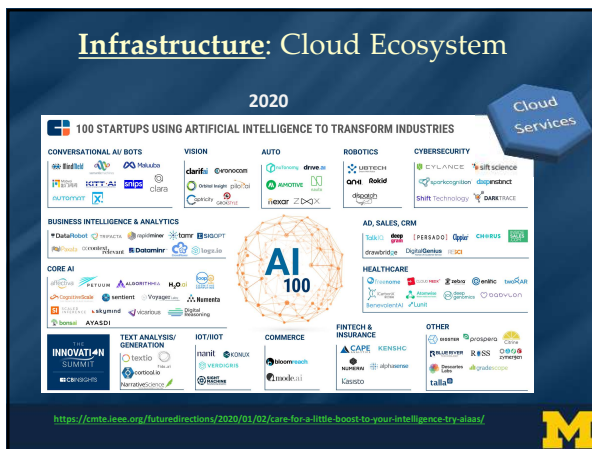
# Why is FAIR Data Sharing Important?

- ❑ 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
- ❑ …

Data

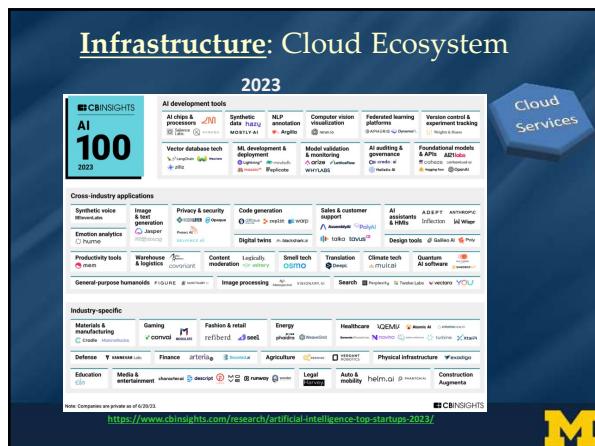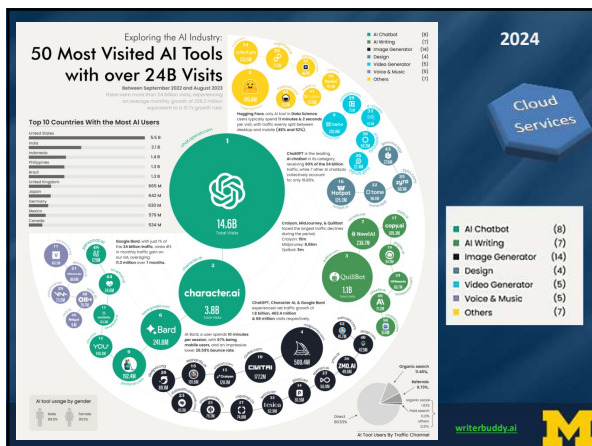FAIR = Findable + Accessible + Interoperable + Reusable

8

## Slide 9

# **Infrastructure**: Cloud Ecosystem

### 2013



Cloud Services

https://socr.umich.edu/docs/BD2K/BigDataResourceome.html

9

## Slide 10

# **Infrastructure**: Cloud Ecosystem

### 2020

100 STARTUPS USING ARTIFICIAL INTELLIGENCE TO TRANSFORM INDUSTRIES



Cloud Services

https://cmte.ieee.org/futuredirections/2020/01/02/care-for-a-little-boost-to-your-intelligence-try-aiaas/

10

## Slide 11

# **Infrastructure**: Cloud Ecosystem

### 2023



Cloud Services

https://www.cbinsights.com/research/artificial-intelligence-top-startups-2023/

11

## Slide 12

Exploring the AI Industry:
## 50 Most Visited AI Tools with over 24B Visits

2024



Cloud Services

| | |
|---|---|
| AI Chatbot | (8) |
| AI Writing | (7) |
| Image Generator | (14) |
| Design | (4) |
| Video Generator | (5) |
| Voice & Music | (5) |
| Others | (7) |

writerbuddy.ai

12

13

## Scholarly Research: OA Pubs/Sharing

- OA Pubs
  - https://en.wikipedia.org/wiki/Open_access
  - https://arxiv.org | https://www.biorxiv.org
  - Blogs (e.g., https://TerryTao.wordpress.com)

- Cloud Services
  - Computing (e.g., Azure, Google, AWS)
  - Storage
  - ICT (information and communication technologies)

- SW
  - https://GitHub.com (e.g., https://github.com/SOCR)
  - http://Cran.r-project.org | Jupyter.org | Rmarkdown.rstudio.com
  - E.g., https://DSPA2.predictive.space

- Licensing
  - https://www.gnu.org/licenses
  - https://socr.umich.edu/html/SOCR_CitingLicense.html

*Pubs*

15

## Rationale for Open Science (Cons)

- Journals impact factor (compared to pay-per-view journals, OA are newer)
- *Predatory* science (dubious quality, profit-centric, spam camouflage)
- Discovery is easy, but validity/utility of the science or tools may be difficult to evaluate *en masse*
- Extra work may be required by all scholars to sift through and identify appropriate materials
- Ambiguity of usage-rights/copyrights/licenses
- Democratization and socialization of science may suffer from some of the same downsides as social-networks
- Is science *competitive* or *collaborative*? Is it a *zero-sum* enterprise?

18

## Rationale for Open Science (Pros)

- We are always stronger together
- Long-term sustainability prefers openness, inclusivity & diversity
- Optimized investments, career advancement, impact & cost-efficiency
- Expeditious discovery, innovation, productization & higher impact
- Rapid devaluation of data-hoarding, clandescent science, knowledge obfuscation
- …



https://www.aaas.org/news/big-data-blog-part-v-interview-dr-ivo-dinov

19

## Rationale for Open Science: Kryder vs. Moore

- Moore's law = the expectation that our computational capabilities, specifically the number of transistors on integrated circuits, doubles approximately every 18-24 months.
- Kryder's law = the volume of data, in terms of disk storage capacity, is doubling every 14-18 months.
- **Kryder ≫ Moore**: Although both laws yield exponential growth, data volume is increasing at a faster pace. Thus, there are clear interests and needs for significant private, public and government engagement in opening, managing, processing, interrogating and interpreting the information content of Big Data.

*Collabs*



Dinov (2016) SMSI | https://www.aaas.org/news/big-data-blog-part-v-interview-dr-ivo-dinov

20

## Reliable, Effective & Secure Data Sharing

- Why is data-sharing difficult?
  monopoly, preservation of *status-quo,* competitive edge, personally identifiable information, IP protection, security (on multiple levels), red tape, …

- FAIR (Findable, Accessible, Interoperable & Reusable) Data are powerful

- Current Data Sharing Landscape?
  Differential Privacy, fully-homomorphic encryption, statistical obfuscation (DataSifter), …

- Digital Twins: https://EHR-Sim.StatisticalComputing.org/clinical-phenotype
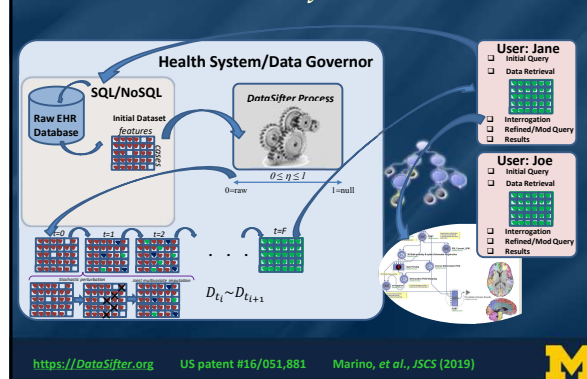
21

## Slide 22

### 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.

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

22

## Slide 23

### DataSifter



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

23

## Slide 24



DataSifter Longitudinal Obfuscator

Try-It:   https://SOCR-DSLO.StatisticalComputing.org

24

## Slide 30

### Case-Studies – General Populations



❑ 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) |

https://www.ukbiobank.ac.uk
https://bd2k.org

30

## Slide 31

### Case-Studies – UK Biobank (Complexities)



Missing Clinical & Phenotypic data for 10K subjects with sMRI, for which we computed 1,500 derived neuroimaging biomarkers.

Including only features observed >30%
(9,914 × 1,475)

Zhou, et al. (2019), SREP   |   https://github.com/SOCR/UKBB_Analytics

31

## Slide 32

### Case-Studies – UK Biobank – NI Biomarkers



32

## Case-Studies – UK Biobank – Successes/Failures

33

## Case-Studies – UK Biobank – Results

*t-SNE plot of the brain neuroimaging biomarkers*

| | | k-means clustering | |
|---|---|---|---|
| | | Cluster 1 | Cluster 2 |
| Hierarchical clustering | Cluster 1 | 3768 (38.0%) | 528 (5.3%) |
| | Cluster 2 | 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 |

34

## 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%) |
| **. . .** | . . . | |
| **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%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%) |

35

## 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.

36

## 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. (2019), SREP*

37

## Complex-Time (*Kime*)

❑ At a given spatial location, $x$, complex time (*kime*) is defined by $\kappa = re^{i\varphi} \in \mathbb{C}$, where:
  ❑ the <u>magnitude</u> represents the longitudinal events order ($r > 0$) and characterizes the longitudinal displacement in time, and
  ❑ event <u>phase</u> ($-\pi \leq \varphi < \pi$) is an angular displacement, or event direction
❑ There are multiple alternative parametrizations of kime in the complex plane
❑ Space-kime manifold is $R^3 \times \mathbb{C}$:
  ❑ $(x, k1)$ and $(x, k4)$ have the same spacetime representation, but different spacekime coordinates,
  ❑ $(x, k1)$ and $(y, k1)$ share the same kime, but represent different spatial locations,
  ❑ $(x, k2)$ and $(x, k3)$ have the same spatial-locations and kime-directions, but appear sequentially in order

38

## Slide 39

### *Kime Parameterizations*



Conjugate Pairs $\{z, \bar{z} \in \mathbb{C}\}$

$z = re^{i\varphi}$
$\bar{z} = re^{-i\varphi}$

$x = (z + \bar{z})/2$
$y = -i(z - \bar{z})/2$
—————
$z = x + iy$
$\bar{z} = x - iy$

$r = \sqrt{z\bar{z}} = \sqrt{|z|z}$

$\varphi = \arccos\frac{z + \bar{z}}{2z\bar{z}}$

$x = r\cos\varphi$
$y = r\sin\varphi$
—————
$r = \sqrt{x^2 + y^2}$
$\varphi = \text{atan2}(y, x)$

Cartesian $\{x, y \in \mathbb{R}^2\}$     Polar $\{(r, \varphi) \in \mathbb{R}^+ \times [-\pi, \pi]\}$

39

## Slide 40

### The Importance of Kime-Magnitude (*time*) and Kime-Phase (*direction*)



**Fourier Analysis**
(real part of the Forward Fourier Transform)

| Square Image Shape | | | | Disk Image Shape | | | |

2D image 1 (square) | Re(FT(square)) | Magnitude FT(square) | Phase FT(square) | 2D image 2 (disc) | Re(FT(disc)) | Magnitude FT(disc) | Phase FT(disc)

**Fourier Synthesis**
(real part of the Inverse Fourier Transform)

Square Image Shape | Disk Image Shape

IFT(FT(square)) ≡ square | IFT using square-magnitude & disc-phase | IFT using square-magnitude & nil-phase | IFT using disc-magnitude & square-phase | IFT using disc-magnitude & nil-phase

40

## Slide 41

### Longitudinal Data Analytics

- **Neuroimaging**:
  - *4D fMRI*: time-series, represents measurements of hydrogen atom densities over a 3D lattice of spatial locations ($1 \leq x, y, z \leq 64$ pixels), about $3 \times 3$ millimeters$^2$ resolution. Data is recorded longitudinally over time ($1 \leq t \leq 180$) in increments of about 3 seconds, then post-processed
  - *State-of-the-art Approaches*: Time-series modeling or Network analysis
  - *Spacekime Analytics*: 5D fMRI kime-series, represent the hydrogen atom densities over the same 3D lattice of spatial locations, longitudinally over the 2D kime space, $\kappa = re^{i\varphi} \in \mathbb{C}$
  - *Differences*: Spacekime analytics estimate and utilize the kime-phases



4D Spacetime     5D Spacekime

4D/5D Reconstructions

Dinov & Velev (2021)

41

## Slide 42

### Spacekime Calculus

- Kime **Wirtinger derivative** (first order kime-derivative at $k = (r, \varphi)$):

In Cartesian coordinates:
$$f'(z) = \frac{\partial f(z)}{\partial z} = \frac{1}{2}\left(\frac{\partial f}{\partial x} - i\frac{\partial f}{\partial y}\right) \quad \text{and} \quad f'(\bar{z}) = \frac{\partial f(\bar{z})}{\partial \bar{z}} = \frac{1}{2}\left(\frac{\partial f}{\partial x} + i\frac{\partial f}{\partial y}\right).$$

In Conjugate-pair basis: $df = \partial f + \bar{\partial} f = \frac{\partial f}{\partial z}dz + \frac{\partial f}{\partial \bar{z}}d\bar{z}$.

In Polar kime coordinates:
$$f'(k) = \frac{\partial f(k)}{\partial k} = \frac{1}{2}\left(\cos\varphi\frac{\partial f}{\partial r} - \frac{1}{r}\sin\varphi\frac{\partial f}{\partial \varphi} - i\left(\sin\varphi\frac{\partial f}{\partial r} + \frac{1}{r}\cos\varphi\frac{\partial f}{\partial \varphi}\right)\right) = \frac{e^{-i\varphi}}{2}\left(\frac{\partial f}{\partial r} - \frac{i}{r}\frac{\partial f}{\partial \varphi}\right)$$

$$f'(\bar{k}) = \frac{\partial f(\bar{k})}{\partial \bar{k}} = \frac{1}{2}\left(\cos\varphi\frac{\partial f}{\partial r} - \frac{1}{r}\sin\varphi\frac{\partial f}{\partial \varphi} + i\left(\sin\varphi\frac{\partial f}{\partial r} + \frac{1}{r}\cos\varphi\frac{\partial f}{\partial \varphi}\right)\right) = \frac{e^{i\varphi}}{2}\left(\frac{\partial f}{\partial r} + \frac{i}{r}\frac{\partial f}{\partial \varphi}\right).$$

- Kime **Wirtinger integration**:

*Path-integral* $\lim_{|z_{i+1} - z_i| \to 0} \sum_{i=1}^{n-1}(f(z_i)(z_{i+1} - z_i)) \cong \oint_{z_a}^{z_b} f(z_i)dz$.

*Definite area integral*: for $\Omega \subseteq \mathbb{C}$, $\int_{\Omega} f(z)dzd\bar{z}$.

*Indefinite integral*: $\int f(z)dzd\bar{z}$, $df = \frac{\partial f}{\partial z}dz + \frac{\partial f}{\partial \bar{z}}d\bar{z}$.

The *Laplacian* in terms of conjugate pair coordinates is $\Delta f = d^2 f = 4\frac{\partial f}{\partial z}\frac{\partial f}{\partial \bar{z}} = 4\frac{\partial f}{\partial \bar{z}}\frac{\partial f}{\partial z}$.

Dinov & Velev (2021)

42

## Slide 43

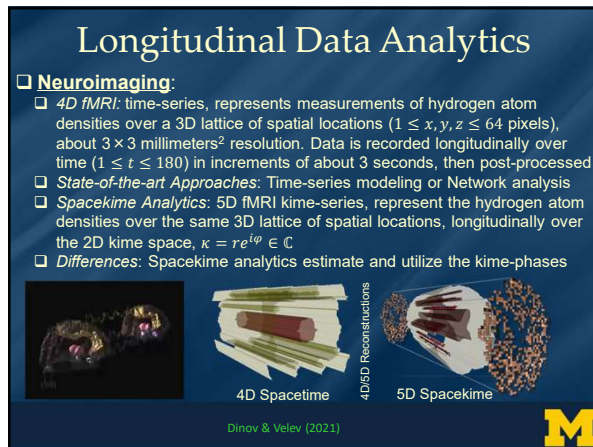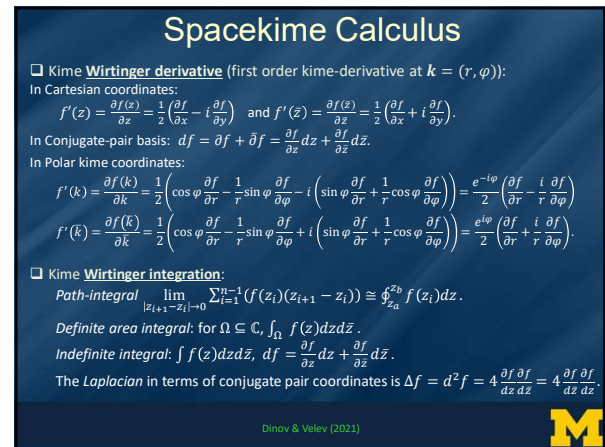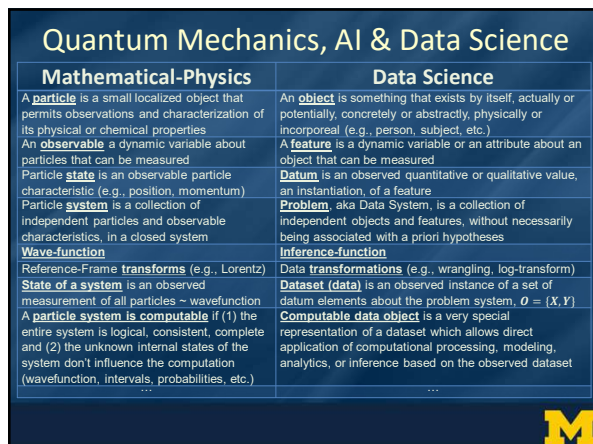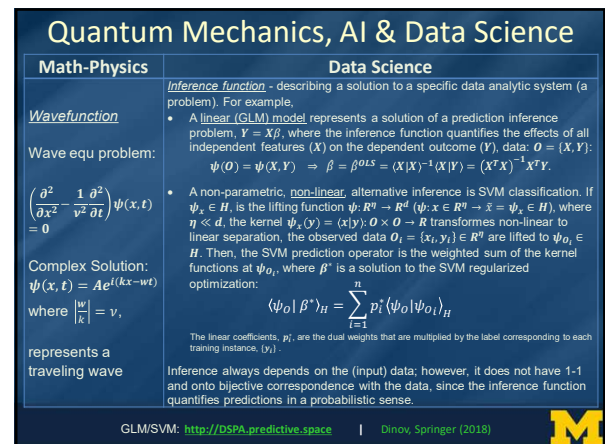### Quantum Mechanics, AI & Data Science

| Mathematical-Physics | Data Science |
|---|---|
| A **particle** is a small localized object that permits observations and characterization of its physical or chemical properties | An **object** is something that exists by itself, actually or potentially, concretely or abstractly, physically or incorporeal (e.g., person, subject, etc.) |
| An **observable** a dynamic variable about particles that can be measured | A **feature** is a dynamic variable or an attribute about an object that can be measured |
| Particle **state** is an observable particle characteristic (e.g., position, momentum) | **Datum** is an observed quantitative or qualitative value, an instantiation, of a feature |
| Particle **system** is a collection of independent particles and observable characteristics, in a closed system | **Problem**, aka Data System, is a collection of independent objects and features, without necessarily being associated with a priori hypotheses |
| **Wave-function** | **Inference-function** |
| Reference-Frame **transforms** (e.g., Lorentz) | Data **transformations** (e.g., wrangling, log-transform) |
| **State of a system** is an observed measurement of all particles ~ wavefunction | **Dataset (data)** is an observed instance of a set of datum elements about the problem system, $O = \{X, Y\}$ |
| 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.) … | **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 … |

43

## Slide 44

### Quantum Mechanics, AI & Data Science

| Math-Physics | Data Science |
|---|---|
| *Wavefunction*<br><br>Wave equ problem:<br><br>$\left(\frac{\partial^2}{\partial x^2} - \frac{1}{v^2}\frac{\partial^2}{\partial t^2}\right)\psi(x,t) = 0$<br><br>Complex Solution:<br>$\psi(x,t) = Ae^{i(kx-wt)}$<br>where $\left|\frac{w}{k}\right| = v$,<br><br>represents a traveling wave | *Inference function* - describing a solution to a specific data analytic system (a problem). For example,<br>• A linear (GLM) model represents a solution of a prediction inference problem, $Y = X\beta$, where the inference function quantifies the effects of all independent features ($X$) on the dependent outcome ($Y$), data: $O = \{X, Y\}$:<br>$\quad \psi(O) = \psi(X,Y) \Rightarrow \beta = \beta^{OLS} = (X|X)^{-1}(X|Y) = (X^TX)^{-1}X^TY.$<br>• A non-parametric, non-linear, alternative inference is SVM classification. If $\psi_x \in H$, is the lifting function $\psi: R^\eta \to R^d$ ($\psi: x \in R^\eta \to \tilde{x} = \psi_x \in H$), where $\eta \ll d$, the kernel $\psi_x(y) = \langle x|y \rangle: O \times O \to R$ transforms non-linear to linear separation, the observed data $O_i = \{x_i, y_i\} \in R^\eta$ are lifted to $\psi_{O_i} \in H$. Then, the SVM prediction operator is the weighted sum of the kernel functions at $\psi_{O_i}$, where $\beta^*$ is a solution to the SVM regularized optimization:<br>$\quad \langle \psi_O \mid \beta^* \rangle_H = \sum_{i=1}^{n} p_i^* \langle \psi_O \mid \psi_{O_i} \rangle_H$<br>The linear coefficients, $p_i^*$, are the dual weights that are multiplied by the label corresponding to each training instance, ($y_i$).<br>Inference always depends on the (input) data; however, it does not have 1-1 and onto bijective correspondence with the data, since the inference function quantifies predictions in a probabilistic sense. |

GLM/SVM: http://DSPA.predictive.space   |   Dinov, Springer (2018)

44

## Slide 45

### *Spacekime* Analytics

- Let's assume that we have:
  (1) Kime extension of Time, and
  (2) Parallels between wavefunctions ↔ inference functions
- Often, we can't directly observe (record) data natively in 5D spacekime.
- Yet, we can measure quite accurately the kime-magnitudes ($r$) as event orders, "times".
- To reconstruct the 2D spatial structure of kime, borrow tricks used by crystallographers [1] to resolve the structure of atomic particles by only observing the magnitudes of the diffraction pattern in k-space. This approach heavily relies on (1) prior information about the kime directional orientation (that may be obtained from using similar datasets and phase-aggregator analytical strategies), or (2) experimental reproducibility by repeated confirmations of the data analytic results using longitudinal datasets.

[1] Rodriguez, Ivanova, Nature 2015

Spacetime → Spacekime Transforms
(1) Phase-estimation
(2) Phase-modeling
(3) Laplace Transform

**5D Spacekime**
3D Space $R^3$
$(x_0, x_1, x_2)$
Observed or Computed

2D Kime $\cong R^2$
$(x_3, x_4)$
Computed

**5D k-space**
3D Space $R^3$
$(f_0, f_1, f_2)$
Observed or Computed

K2 Kaluza-Klein $\cong R^2$
$(time\ (t), phase\ (\phi))$
observed directly   estimated

Data Science Analytics — Experimental Science

45

## Slide 46

### *Spacekime* Analytics: fMRI Example

- 3D isosurface Reconstruction of (2D space, 1D time) fMRI signal



**4D spacetime:** Reconstruction using trivial phase-angle; kime=time=(magnitude, 0)

**5D Spacekime:** Reconstruction using correct kime=(magnitude, phase)

3D pseudo-spacetime reconstruction:

$$f = \hat{h}\ (\underbrace{x_1, x_2}_{space},\ \underbrace{t}_{time})$$

46

## Slide 47

### *Spacekime* Analytics: Kime-series = Surfaces (not curves)

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

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



47

## Slide 48

### Bayesian Inference Representation

- We can formulate spacekime inference as a Bayesian parameter estimation problem:

$$\underbrace{p(\gamma|X, \varphi')}_{posterior\ distribution} = \frac{p(\gamma, X, \varphi')}{p(X, \varphi')} = \frac{p(X|\gamma, \varphi') \times p(\gamma, \varphi')}{p(X, \varphi')} = \frac{p(X|\gamma, \varphi') \times p(\gamma, \varphi')}{p(X|\varphi') \times p(\varphi')} =$$
$$= \frac{p(X|\gamma, \varphi')}{p(X|\varphi')} \times \frac{p(\gamma, \varphi')}{p(\varphi')} = \frac{p(X|\gamma, \varphi') \times p(\gamma|\varphi')}{\underbrace{p(X|\varphi')}_{observed\ evidence}} \propto \underbrace{p(X|\gamma, \varphi')}_{likelihood} \times \underbrace{p(\gamma|\varphi')}_{prior}.$$

- In Bayesian terms, the posterior probability distribution of the unknown parameter $\gamma$ is proportional to the product of the likelihood and the prior.

- In probability terms, the posterior = likelihood times prior, divided by the observed evidence, in this case, a single spacetime data point, $x_{i_o}$.

48

## Slide 49

### Spacekime Analytics using fMRI

- Complex-valued *finger tapping* fMRI ($64x\ 64y\ 40z\ 160t$)

fMRI time-series forecasting

Temporal Dynamics of a Voxel in Somatosensory Motor Area



On-Off fMRI time-series to Kimesurface

49

## Slide 50

### Try Spacekime Analytics using Sim Data



https://Kime.StatisticalComputing.org

50

## Foundational Generative AI Models (GAIMs)



https://socr.umich.edu/GAIM/

52

## SOCR Augmented Intelligence Agent



https://AIA3.statisticalcomputing.org/
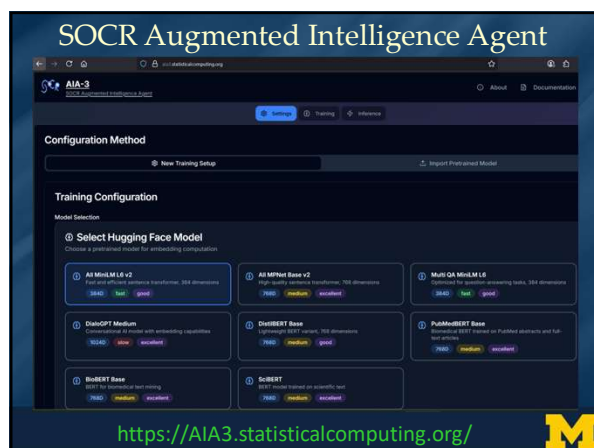
53

## What's Next?

- Lots of "open problems" in data-science, e.g., fundamentals of data representation & analytics
- The SOCR team is developing:
  - Compressive Big Data Analytics (CBDA) technique – an ensemble learning meta-algorithm
  - DS Time-Complexity and Inferential-Uncertainty
- Need lots of community, institutional, state, federal, and philanthropic support to advance <u>open data science</u> methods, enhance the computing infrastructure, train/support students/fellows, and tackle the *Kryder Law ≫ Moore Law* trend
- **Web**:        **www.SOCR.umich.edu**
- **Git**:        **https://github.com/SOCR**
- **Projects:** **www.socr.umich.edu/html/SOCR_Research.html**
- **Apps:**      **https://socr.umich.edu/HTML5/**

55

## Acknowledgments

Slides Online: "SOCR News"

https://SOCR.umich.edu

56