

Trustworthy AI through veridical data science and interpretable machine learning

Bin Yu Statistics and EECS, UC Berkeley

Virtual Meeting on Statistics and AI Statistics Committee of the Hungarian Academy of Science Sept. 16, 2021

1

Al is part of modern life

Virtual assistants	Wearable	health devices	Recommendation systems
			(Tourube, Facebook)
Online news	Bill Gates: A.I. is i 'both promising a Published Two, May 26 2019-8145 AM EDT - Updated Two, May 26	ike nuclear energy — Ind dangerous'	Election campaigns
Self-driving cars	Catherine Clifford @CATCLIFFORD	Shore 🗲 🐭 in 😂	Online gaming
Precision medicine Chemistry			Biology Neuroscience
Materials S	science	Law	Sociology
Cosmology	Economics	Political Science	and beyond

Data science (DS) is a key element of AI



Conway's Venn Diagram

Goal:

Leverage algorithms to combine data with domain knowledge to make decisions and generate new knowledge

Trustworthy AI: two complementary approaches

- "Best practices" to maximize the promise (precaution)
- Damage control to **reduce** the **danger** (intervention)

Veridical Data Science (Y. and Kumbier, PNAS, 2020)

Extracts reliable and reproducible information from data, with an enriched technical language to communicate and evaluate empirical evidence in context of human decisions and domain knowledge

Necessary to realize the promises and mitigates the dangers of AI

Frontiers in biological and medical DS

T0955 / 5W9F





medium.com

T0954 / 6CVZ

Structures: Ground truth (green) Predicted (blue)





T0965 / 6D2V

https://deepmind.com/blog/alphafold/

Machine Learning and Personalization



Website of S. Saria at JHU

Scientific ML (sML): part of trustworthy AI

- It uses machine learning/statistics for scientific research to extract, from data, discoveries, theory, and knowledge
- It builds scientific principles/theory in machine learning algorithms
- It iterates between the above two steps
- Results are subject to scientific standards



Multi-scale deep learning and single-cell models of cardiovascular health

PIs: Euan Ashley, Rima Arnaout, Ben Brown, Atul Butte, James Priest, Bin Yu Collaborators: Chris Re, Deepak Srivastava





M. Behr



K. Kumbier







N. Youlton









M. Aguirre

Palomera









C. Weldy

A. Agarwal







O. Ronen



8

W. Hughes

Biohub project as a data science problem

- Medical question: which genes interact to induce HCM (hypertrophic cardiomyopathy)?
- Which **data** to use? How to clean?
- EDA: summaries, plots, ...
- **Modeling:** Which algorithms to use to find nonlinear interactions?
- Interpretation & evaluation of results



Data Science Life Cycle (DSLC): A holistic view



To maximize promise,

Scientific conclusions must **reflect reality** and be **stable** to human judgment calls throughout the integrated data science life cycle (DSLC).

Thus, data science requires **quality control** and **standardization** inspired by empirical practice.

2001

Statistical Science 2001, Vol. 16, No. 3, 199–231



Statistical Modeling: The Two Cultures

Leo Breiman

The Algorithmic Modeling Culture

The analysis in this culture considers the inside of the box complex and unknown. Their approach is to find a function $f(\mathbf{x})$ —an algorithm that operates on \mathbf{x} to predict the responses \mathbf{y} . Their black box looks like this:

The Data Modeling Culture

The analysis in this culture starts with assuming a stochastic data model for the inside of the black box. For example, a common data model is that data are generated by independent draws from

response variables = f(predictor variables, random noise, parameters)

Machine learning

Statistics

Cross-validation

Deep Learning, AlphaGo, AlphaFold, self-driving cars, ... Linear model, Logistic regression, PCA, p-value, t-test, ...

Predictability, Computability, Stability (PCS) framework

one culture

for veridical data science

Rest of the talk

PCS conceptual framework:

veridical data science with reliable and reproducible results

Interpretable Machine Learning

What is interpretable ML or iML? iterative random forests (iRF) (based on PCS) Adaptive wavelet distillation (AWD) of deep learning networks

PCS-based software:

Veridical Flow and simChef

2001

Statistical Science 2001, Vol. 16, No. 3, 199–231



Statistical Modeling: The Two Cultures

Leo Breiman

The Algorithmic Modeling Culture

The analysis in this culture considers the inside of the box complex and unknown. Their approach is to find a function $f(\mathbf{x})$ —an algorithm that operates on \mathbf{x} to predict the responses \mathbf{y} . Their black box looks like this:

The Data Modeling Culture

The analysis in this culture starts with assuming a stochastic data model for the inside of the black box. For example, a common data model is that data are generated by independent draws from

$$\label{eq:response variables} \begin{split} \text{response variables} = f(\text{predictor variables},\\ \text{random noise, parameters}) \end{split}$$

Machine learning



Statistics

Deep Learning, AlphaGo, AlphaFold, self-driving cars, ... Linear model, Logistic regression, PCA, p-value, t-test, ...

Image credit: https://www.lib.berkeley.edu/news_events/bridge/sfobay.html

PCS framework for veridical data science: One culture Y. and Kumbier (PNAS, 2020)



Three principles of data science:

(P)redictability [ML and Stats]

(C)omputability [ML]

(S)tability [Stats]

Unifies, streamlines, and expands ideas and best practices in **both** ML and Stats

Veridical Data Science



Image credit: R. Barter

The stability principle

Reproducibility is imperative for any scientific discovery. More often than not, modern scientific findings rely on statistical analysis of high-dimensional data. At a minimum, reproducibility manifests itself in **stability** of statistical results relative to **reasonable perturbations** to data and to the model used.

- Yu (2013) [Stability]

PCS connects science with engineering

• **Predictability** and **stability** embed the scientific principles of *prediction* and *replication*





Computability is a practical necessity and includes data-inspired simulations



Image credits: nstat.org, hub.jhu.edu, vox.com, Andras Libal

PCS in a nutshell

Predictability for reality checkStability analysis over human judgment callsComputability is implicit in P and S

Data Science Life Cycle



Image credits: R. Barter and toronto4kids.com

PCS in a nutshell

Predictability for reality checkStability analysis over human judgment callsComputability is implicit in P and S

Stability analysis "shakes" every part of DSLC to check robustness

"Stability analysis" defined **broadly** and **flexibly**

Data Science Life Cycle



Image credits: R. Barter and toronto4kids.com

PCS documentation [on GitHub (JupyterNotebook)] bridges reality and models



Image credits: Rebecca Barter

What is PCS to a doctor?

"The PCS framework builds a **working relationship** between data and the clinical world."

"PCS is a 'look under the hood' to ensure that the **conclusions found** are what **the data genuinely suggest**. In all, PCS is a holistic approach to helping the clinician **understand**, **interpret**, **and build the science** we need to help our patients."

Dr. Aaron Kornblith, ER, UCSF



main medical collaborator on PCS stress-testing of PECARN CDR



Data preprocessing (shoe-leather work) causes data perturbations

American Economic Review: Papers & Proceedings 100 (May 2010): 573–578 http://www.aeaweb.org/articles.php?doi=10.1257/aer.100.2.573

Growth in a Time of Debt

By CARMEN M. REINHART AND KENNETH S. ROGOFF*

RR paper was covered widely in popular media, often as "high debt/GDP ratio is bad for growth," to support austerity policies in UK and Europe.

Herdon, Ash and Pollin (2014) found that RR had exclusive data selection (cleaning), coding errors, and unconventional weighting.

When corrected, RR's conclusion fails to hold.

Data choice causes data perturbations



Low resolution



HCM cell size study after knockout experiments by Ashley Lab at Stanford

Data partitioning causes data perturbations

different people create data splits in different ways

Data perturbations to assess P

Data splitting (stratified by treatment & outcome)

Training folds (cross-validation) Works well if data units are symmetric



Dwivedi et al (2020) Stable discovery of interpretable subgroups via calibration in causal studies. ²⁶ International Statistical Review

Choice of method/algorithm causes model perturbations:

different people prefer different methods/algorithms

Some recommendations from PCS

- Make sure "model" means the same thing for a multi-disciplinary team
- Keep multiple versions of annotated or cleaned data
- Keep multiple models after some prediction/reality check, such as in climate modeling



The change in global-mean temperature estimated by nine climate models forced by the SRES A2 emission scenario. (Source: IPCC TAR, Chapter 9)

Expanding statistical inference under PCS

Modern goal of **statistical inference** is to provide one source of **evidence** to domain experts **for decision-making.**

The key is to provide data evidence in a **transparent** manner so that **domain experts can understand** as much as possible the data evidence generation **to evaluate the strength** of evidence.

Traditionally, the **p-value** has been used as evidence for decisions, but its use has been so **problematic** that psychology journals have **banned** it.

PCS Inference (Yu and Kumbier, 2020)

Predictability

Use prediction error for model checking

Stability

Assessed across data and model perturbations (with data perturbation broadly interpreted)

Computability

Implicitly required by P and S

Includes data-inspired simulation

Does not assume a probabilistic generative model -- similar to bootstrap-based inference when the probabilistic model approximates reality well

Interpretable ML (iML): necessary for sML and trustworthy AI

(Murdoch, Singh, Kumbier, Abbasi-Asl, and Y., PNAS, 2019) "Definitions, Methods and Applications in Interpretable Machine Learning"



"We define interpretable machine learning as the extraction of relevant knowledge from a machine-learning model concerning relationships either contained in data or learned by the model. Here, we view knowledge as being relevant if it provides insight for a particular **audience** into a **chosen problem**. These insights are often used to guide communication, actions, and discovery."

iML through PDR: **Predictive accuracy**, **Descriptive accuracy and Relevancy**

iML-PDR in one figure



Two vignettes: from predictive models to interpretations

Q: what genes drive a heart disease HCM?

Iterative random forests to discover predictive and stable high-order interactions

Sumanta Basu^{a,b,c,1}, Karl Kumbier^{d,1}, James B. Brown^{c,d,e,f,2}, and Bin Yu^{c,d,g,2}

Co-authors





S. Basu

K. Kumbier



B. Brown



Culmination of 3+ years of work

Pattern Recognition vs. Pattern Discovery

Pattern Recognition: Finding something for which you already know to look



Pattern Discovery: Identifying structure that hasn't been seen before



Order-4 interaction regulates *eve* stripe 2



Goto et al. (1989), Harding et al. (1989), Small et al. (1992), Isley et al. (2013), Levine et al. (2013)

Capturing the form of genomic interactions

- Interactions are high-order and combinatorial in nature
- Interactions can vary across space and time as biomolecules carry out different roles in varied contexts
- Morphogen Conc. (M) Morphogen Conc. (M)

(Wolpert, 1969; Jaeger and Reinitz, 2006)



(Hartenstein, 1993)

 Interactions exhibit thresholding behavior, requiring sufficient levels of constitutive elements before activating



From genomic to statistical interactions





 ${\rm Order-}_S \ {\rm interaction,}$

 $S \subseteq \{1, \dots, p\}, |S| = s$

Random Forests (RF) (Breiman, 2001)

Draw *T* bootstrap samples and fit a modified CART to each sample.

- 1. Grow CART trees to purity.
- When selecting splitting feature, choose a subset of mtry features uniformly at random and optimize CART criterion over subsampled features.



Iterative random forests (iRF) finds predictive and stable interactions for a Drosophila enhancer prediction problem



80% of pairwise gene interactions are validated by past biological experiments in the literature Three-way interactions are recommendations

Iterative random forests (iRF) keep predictive accuracy, and finding stable interactions for a fruitfly enhancer prediction problem from TFs



80% of pairwise interactions are validated by past biological experiments in the literature Three-way interactions are recommendations siRF-estimated fruitfly TF binding will be made available as UCSC genome browser track (Kumbier et al, 2021)



Genome Browser

Biohub project for cardiovascular health

Berkeley

CHAN ZUCKERBERG

Cardiovascular phenotype LVM from cardiac MRI: (n = 30,000 in UKBiobank)

Stanford

- Left Ventricle Mass (LVM): proxy for a well-known heart disease called Hypertrophic Cardiomyopathy (HCM) (rare variant with 1 in 500)
- Genetic association complex: predictability and stability both low

Yu Group at Berkeley found 4 **predictive and stable gene pairs** that might drive LVM, discovered using their method of iterative random forests (iRF).

Ashley Lab at Stanford Medical School are carrying out **siRNA transfection experiments** with promising preliminary results.



Adaptive wavelet distillation (AWD) from neural networks through interpretations

Co-authors



W. Ha,

C. Singh,

F. Lanusse,

S. Upadhyayula

https://arxiv.org/abs/2107.09145

CODE: https://github.com/Yu-Group/adaptive-wavelet-distillation

AWD goal

Given a domain where a Deep Neural Network predicts well, distill it into a simple **learned** wavelet transform

Improves interpretability, compression, and efficiency

Adaptive Wavelet Distillation (AWD)



AWD finds h and g to minimize L(h, g) through (stochastic) gradient descent

Validation of AWD methodology in a biology problem

Clathrin-mediated endocytosis (CME)



CME ``is a key process in vesicular trafficking that transports a wide range of cargo molecules from the cell surface to the interior.''

-- Kaksonen and Roux (2018) *Nature Reviews*. https://www.nature.com/articles/nrm.2017.132

Tracking molecular partners is a central problem in cell biology

...but is experimentally difficult

we aim to predict molecular partners\ (from clathrin to auxilin)



AWD performs best: prediction, compression, speed-up



R² score

For each scale, we select 6 largest coefficients or 30 coefficients

- proportion of variance in Y that is explained by the model

AWD	DB5 Wavelet	LSTM	AWD w/o interp. loss
0.263	0.197	0.237	0.231

A large build up in clathrin fluorescence followed by sharp drop is a highly predictive signature of a successful CME event

PCS Softwares: easy to use

Design Principles:

Transparent (P) Realistic (P)

Intuitive (**C**) Modular (**C**) Efficient (**C**)

Reproducible (S)





Veridical Flow (on-going): a PCS wrapper for ML in Python

Make prediction screening and stability analysis **simple**, **reproducible**, and **computationally efficient**

Create intuitive capabilities to **inspect**, **manipulate**, and **visualize** the ML pipeline



github.com/Yu-Group/veridical-flow

Veridical Flow: a PCS wrapper for ML in Python

Agnostic to underlying ML framework

Closely integrated with capabilities of Ray & MLflow (both from CS at Berkeley)

- Distributed computation 🔗 RAY
- Tracking results across perturbations
- Visualization of performance metrics
- Packaging trained models for reproducibility and real-world deployment



github.com/Yu-Group/veridical-flow

simChef (on-going): an R package for PrinCipled Simulations



- Provide **powerful tools** for **evaluating methods** efficiently across a **variety of scenarios** (perturbations)
- Encourage adherence to principles of strong, reliable, trustworthy data science



github.com/Yu-Group/simChef

simChef Interactive PCS Documentation

• Automated generation of transparent and reproducible documentation in R Markdown for rapid and veridical communication of scientific results.

Regression Experiment	Simulation Experiment Recipe
Simulation	Objectives Data Generation Methods and Evaluation Visualizations
Recipe	The objective of this simulation experiment is to provide a toy example on how to use surched and showcase the automated R Markdown- generated documentation. For the sake of illustration, this toy simulation experiment studies the performance of linear regression at the
Base Linear Regression Experiment	surface-level with the sole purpose of facilitating an easy-to-understand waikthrough. [Typically, the objective of the simulation experiment (and this blurb) will be more scientific than instructive and will warrant additional context/background and domain knowledge.]
Linear Gaussian DGP	

PCS inference case study: epiTree for epistasis discovery



Post-docs & Students:









M. Behr

K. Kumbier M. Aguirre

A. Cordova-Palomera

Learning epistatic polygenic phenotypes with Boolean interactions

https://www.biorxiv.org/content/10.1101/2020.11.24.396846v1

epiTree uses iRF for interaction selection and new PCS p-value for epistasis testing

Summary

- Veridical data science is necessary for trustworthy AI
- **PCS framework** (with **documentation on GitHub**) for Veridical data science (8 successful PCS case studies including iRF)
- Scientific ML is part of trustworthy AI: it needs interpretable ML
- iML based on PDR, decision trees from iRF, AWD from DNNs
- **Domain knowledge** is important, and building on this, **PCS** generates testable results for **external validation** (e.g., experiments and follow-up studies)
- Importance of reliable and easy-to-use **software**
- Trustworthy AI needs a fair reward system and incentives

Software from the Yu Group stat.berkeley.edu/~yugroup/code.html

- iRF, siRF
- epiTree
- Superheat
- Veridical Flow (on-going)
- simChef: PCS Simulation package in R (on-going)

Also, newly accepted by the Journal of Open Source Software (with 10k+ downloads, and 240+ github stars):

imodels (Python) for concise, transparent, and accurate predictive modeling: provides a simple interface for fitting & using state-of-the-art interpretable models compatible with scikit-learn.

Hope PCS is useful for your projects

Bin Yu website: <u>https://binyu.stat.berkeley.edu/</u> Papers: <u>https://binyu.stat.berkeley.edu/papers</u>

- 1. B. Yu and K. Kumbier (2020), **"Veridical data science"**, PNAS. --- PCS framework
- 2. S. Basu, K. Kumbier, B. Brown and B. Yu (2018). **"Iterative random forests to discover predictive and stable high-order interactions",** PNAS

3. K. Kumbier, S. Basu, J. Brown, S. Celniker, B. Yu (2018) Refining interaction search through signed iterative Random Forests (signed iRF or siRF) <u>https://arxiv.org/abs/1810.07287</u>

4. M. Behr, K. Kumbier, M. Aguirre, R. Arnaout, E. Ashley, A. Butte, R. Arnout,
B. Brown, J. Priest, B. Yu (2020). "Learning epistatic polygenic phenotypes with
Boolean interactions" <u>https://www.biorxiv.org/content/10.1101/2020.11.24.396846v1</u>

5. Singh, Ha and Yu (2021). Interpreting and improving deep-learning models with reality checks (iML review paper including AWD) <u>https://arxiv.org/abs/2108.06847</u>

Book using PCS for DSLC by Yu and Barter with MIT Press free online interactive copy (plan: 2021 fall)

Veridical Data Science: A Book

Bin Yu¹³ and Rebecca Barter¹

Department of Statistics, UC Berkeley ²Department of Electrical Engineering and Computer Science, UC Berkeley





What skills does the book teach?

Vericical Data Science (VDS) will teach the critical thinking, analytic, human-interaction and communication skills required to effectively formulate protriems and find reliable and trustworthy solutions. VDS explains concepts using visuals and plain English, rather than math and code.

The primary skills taught are:



larit will incore the

Formulate anowerstore quastions using the para available beindorrow all analytic concerns and rescale Discionant all adultate decisions -Approximate common textremans to universitive setunteers Deut with steal, meany utata



Communication



Exploratory Visual Summarian Property and in many worst and Appropriate Including the Automotion of tals and finitings to an enternal

ARCTINESS TAXABLE Recording action assigns records for these excellent. Dated St Last, Pauly 1984

Core guiding principles for the book Three realms

The DS Lifecycle



The Data Science Lifecycle is an iterative process that takes the analyst from problem formulation, data cleaning, exploration, algorithmic analysis, and finally to obtaining a verifiable solution that can be used for future decision-making.

Bending together concepts from statistics, computer science and domain knowledge. the data science life cycle is an terative process that involves human analysis learning from data and refining their project-specific questions and snalytic approach as they learn.

Intended Reader/Audience

Anyone who wants to learn the intuition and critical thinking skills to become a data scientist or work with data scientists.

Nother a mathematical nor a coding background is required. VD5 could form the basis of a semester- or multi-semester long introductory data science university course, either as an upper division undergraduate or early graduatelevel course.

Interested? Get in touch!

Bin Yu

Website: https://www.stat.toerkeley.edu#binyu/Site/Welcome.html



Readers will learn to view every data problem through the lens of connecting the three sealing:

- d) the question being asked and the data collected (and the reality the data representab
- (2) the algorithms used to represent the dota
- (3) future data on which these algorithms will be used to guide decision-making. Computability accorthenic and data Guiding the reader to connect the three realms is a means of guiding the reader through the data science lifecycle.



The PCS framework provides concrete techniques for finding evidence for the connections between the three realms. Predictability: If the patterns found in the original data also appear in withheld or new data, they are said to be predictable. if an analysis or algorithm finds predictable patterns, then these patterns are likely to be capturing real phenomena. efficiency and scalability is essential to ensuring that the results and solutions is g. a predictive algorithm) can be efficiently applied to new data. Stability: minimum requirement for reproducibility. If results change in the presence of minor modifications of the della (e.g. via perturbations) or human analytic decisions, then there might not be a strong connection between the analysis/aporithms and the reality that underlies the data.

Rebecca Barter

Email rebecceberter@berkele_59k Website: www.selbeccalterter.com Twitter: Grittanier

Email: breyu@stat.berkeley.edu

Berkeley Computing, Data Science, and Society (CDSS led by Jennifer Chayes)

Academics * Research * News * Events * About * Support Us





A data science program for everyone Data8 on EdX

Berkele

Professional Certificate in Foundations of Data Science The Future of Data Associate Provost Jennifer Chayes on equality, equity, opportunity in data

Data Science Major Data8 (1000+ students) Data100 (1000+) Data102 (200+) co-created and co-taught by Stats and EECS faculty

I'm interested O

What you will learn

- · How to interpret and communicate data and results using a vast array of real-world examples from different domains
- How to make predictions using machine learning and statistical methods.
- Computational thinking and skills, including the Python programming language for analyzing and visualizing data
- How to think critically about data and draw robust conclusions based on incomplete information





Self-paced Progress at your own speed



4 months 4 - 6 hours per week



\$537.30 \$597 USD For the full program experience