

What is health data science?

STA 198: Introduction to Health Data Science

Yue Jiang

May 18, 2020

The following material was used by Yue Jiang during a live lecture.

Without the accompanying oral comments, the text is incomplete as a record of the presentation.

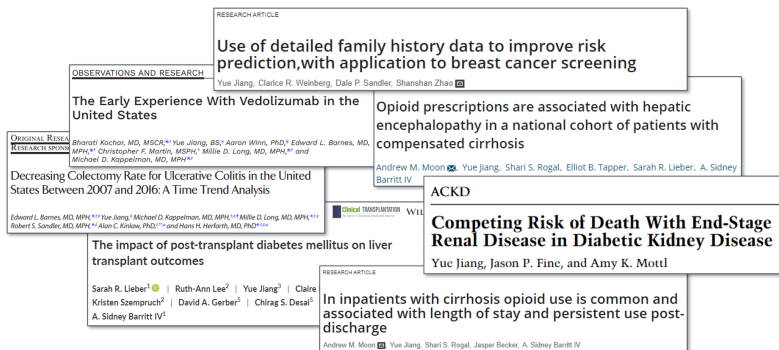
Who are you?



Who are you?



Who are you?



What have I gotten myself into?

STA 198L is a rigorous ten-week online introduction to health data science that...

- ▶ ...provides a tour of basic statistical methods commonly encountered in public health and biomedical research,
- ▶ ...emphasizes understanding of the methods, using them to arrive at data-driven decisions, effective communication of such results, and critically assessing existing evidence,
- ▶ ...is motivated by timely, relevant examples from current research in biomedicine, epidemiology, and public health policy, and
- ▶ ...utilizes modern software such as RStudio and GitHub to reproducibly examine and manipulate data to make sound scientific conclusions.

What IS health data science?

A process that converts data into useful information, whereby practitioners

1. form a question of interest,
2. collect and summarize data,
3. and interpret the results

RHEUMATISM POSITIVELY CURED,
Also Gout, Sciatica, Neuralgia, Numbness, and Blood Disorders, resulting from excesses, impaired circulation, or sluggish liver, by wearing the genuine **Dr. BRIDGMAN'S**


full-power **ELECTRO-MAGNETIC RING**, a quick and reliable remedy, as thousands testify, and it **WILL CURE YOU.**

"Offices of the New York Bottling Co., N.Y."
"Dr. Bridgman's Ring quickly cured me after years of intense suffering from Rheumatism. Ten thousand dollars would not buy mine if I could not obtain another. I confidently recommend it to all who have Rheumatism."
"GEO. W. RAYNER, PRES."
"Dr. Bridgman's Ring has performed most miraculous cures of Rheumatism and Gout."
"O. VANDER BILT, N.Y."
"I have not had a twinge of Rheumatic Gout since wearing Dr. Bridgman's Ring. It is a quick cure."
"JUDGE REYNOLDS, N.Y. CITY."
Thousands of others offer similar testimony.

We have supplied these rings to *Harrison, Cleveland, Blaine, Depew, Bismarck*, and other eminent men. Their effect is marvellous. Price, \$1.00 plain finish, and \$2.50 heavy gold-plated. All sizes. For sale by **Druggists and Jewelers**, or we will mail it, post-paid, on receipt of price and size.

There is absolutely no other ring but **Dr. Bridgman's** possessing real merit for the cure of Rheumatism. Beware of Imitations.

THE A. BRIDGMAN CO. {373 Broadway, N. Y., and
{1224 Masonic Temple, Chicago.



What is health data science good for?

The NEW ENGLAND JOURNAL of MEDICINE

SPECIAL ARTICLE

Mortality in Puerto Rico after Hurricane Maria

Nishant Kishore, M.P.H., Domingo Marqués, Ph.D., Ayesha Mahmud, Ph.D.,
Mathew V. Kiang, M.P.H., Irmay Rodriguez, B.A., Arlan Fuller, J.D., M.A.,
Peggy Ebner, B.A., Cecilia Sorensen, M.D., Fabio Racy, M.D., Jay Lemery, M.D.,
Leslie Maas, M.H.S., Jennifer Leaning, M.D., S.M.H., Rafael A. Irizarry, Ph.D.,
Satchit Balsari, M.D., M.P.H., and Caroline O. Buckee, D.Phil.

NEJM (July, 2018)

What is health data science good for?

ARTICLES

<https://doi.org/10.1038/s41477-018-0263-1>

nature
plants

Decreases in global beer supply due to extreme drought and heat

Wei Xie ^{1*}, Wei Xiong^{2,3,4}, Jie Pan ², Tariq Ali¹, Qi Cui⁵, Dabo Guan ^{6,7*}, Jing Meng ⁸,
Nathaniel D. Mueller⁹, Erda Lin ^{2*} and Steven J. Davis^{9,10}

Nature *Plants* (October, 2018)

What is health data science good for?

SCIENCE TRANSLATIONAL MEDICINE | RESEARCH ARTICLE

GENETIC DIAGNOSIS

Diagnosis of genetic diseases in seriously ill children by rapid whole-genome sequencing and automated phenotyping and interpretation

Michelle M. Clark¹, Amber Hildreth^{1,2,3}, Sergey Batalov¹, Yan Ding¹, Shimul Chowdhury¹, Kelly Watkins¹, Katarzyna Ellsworth¹, Brandon Camp¹, Cyrielle I. Kint⁴, Calum Yacoubian⁵, Lauge Farnaes^{1,2}, Matthew N. Bainbridge^{1,6}, Curtis Beebe⁷, Joshua J. A. Braun¹, Margaret Bray⁸, Jeanne Carroll^{1,2}, Julie A. Cakici¹, Sara A. Caylor¹, Christina Clarke¹, Mitchell P. Creed⁹, Jennifer Friedman^{1,10}, Alison Frith⁵, Richard Gain⁵, Mary Gaughran¹, Shauna George⁷, Sheldon Gilmer⁷, Joseph Gleeson^{1,10}, Jeremy Gore¹, Haiying Grunenwald¹², Raymond L. Hovey¹, Marie L. Janes¹, Kejia Lin⁷, Paul D. McDonagh⁸, Kyle McBride⁷, Patrick Mulrooney¹, Shareef Nahas¹, Daeheon Oh¹, Albert Oriol⁷, Laura Puckett¹, Zia Rady¹, Martin G. Reese¹³, Julie Ryu^{1,2}, Lisa Salz¹, Erica Sanford^{1,2}, Lawrence Stewart⁷, Nathaly Sweeney^{1,2}, Mari Tokita¹, Luca Van Der Kraan¹, Sarah White¹, Kristen Wigby^{1,2}, Brett Williams⁵, Terence Wong¹, Meredith S. Wright¹, Catherine Yamada¹, Peter Schols⁵, John Reyniers⁸, Kevin Hall¹², David Dimmock¹, Narayanan Veeraraghavan¹, Thomas Defay⁸, Stephen F. Kingsmore^{1*}

Copyright © 2019
The Authors, some
rights reserved;
exclusive licensee
American Association
for the Advancement
of Science. No claim
to original U.S.
Government Works

Science Translational Medicine (April, 2019)

What is health data science good for?



Annals of Oncology 29: 1836–1842, 2018
doi:10.1093/annonc/mdy166
Published online 28 May 2018

ORIGINAL ARTICLE

Man against machine: diagnostic performance of a deep learning convolutional neural network for dermoscopic melanoma recognition in comparison to 58 dermatologists

H. A. Haenssle^{1*,†}, C. Fink^{1†}, R. Schneiderbauer¹, F. Toberer¹, T. Buhl², A. Blum³, A. Kalloo⁴,
A. Ben Hadj Hassen⁵, L. Thomas⁶, A. Enk¹ & L. Uhlmann⁷

Annals of Oncology (August, 2018)

What is health data science good for?

Ramucirumab plus erlotinib in patients with untreated, EGFR-mutated, advanced non-small-cell lung cancer (RELAY): a randomised, double-blind, placebo-controlled, phase 3 trial



Kazuhiko Nakagawa, Edward B Garon, Takashi Seto, Makoto Nishio, Santiago Ponce Aix, Luis Paz-Ares, Chao-Hua Chiu, Keunchil Park, Silvia Novello, Ernest Nadal, Fumio Imamura, Kiyotaka Yoh, Jin-Yuan Shih, Kwok Hung Au, Denis Moro-Sibilot, Sotaro Enatsu, Annamaria Zimmermann, Bente Fridmødt-Møller, Carla Visseren-Gruel, Martin Reck, for the RELAY Study Investigators*

Lancet *Oncology* (October, 2019)

What is health data science good for?

Research

JAMA Internal Medicine | [Original Investigation](#)

Comparison of Hospital Mortality and Readmission Rates for Medicare Patients Treated by Male vs Female Physicians

Yusuke Tsugawa, MD, MPH, PhD; Anupam B. Jena, MD, PhD; Jose F. Figueroa, MD, MPH; E. John Orav, PhD; Daniel M. Blumenthal, MD, MBA; Ashish K. Jha, MD, MPH

JAMA Internal Medicine (February, 2017)

Back to the statistical process

1. Forming a question of interest
2. Collecting and summarizing data
3. Interpreting the results



Identifying the population and question of interest

The **population** is the group we'd like to learn something about:

- ▶ What is the prevalence of diabetes among **U.S. adults**, and has it changed over time?
- ▶ Is there a relationship between tumor type and five-year mortality in **breast cancer patients**?
- ▶ Does the average amount of caffeine vary by vendor in **12 oz. cups of coffee at Duke coffee shops**?

If we had data from every unit in the population, we could just calculate what we wanted and be done!

Sampling from the population

Unfortunately, we (usually) have to settle with a **sample** from the population.

Ideally, the sample is **representative**, allowing us to use **probability and statistical inference** to make conclusions that are **generalizable** to the broader population of interest.

Sampling methods

Probability sampling (e.g., simple random sampling, stratified, cluster, or multi-stage sampling)

- ▶ All units have a known chance of being selected
- ▶ More likely to be generalizable
- ▶ Can be more expensive and time-consuming

Non-probability sampling (e.g., quota, convenience, or snowball sampling)

- ▶ Some units unable to be selected, with no way of knowing size or effect of sampling errors
- ▶ Less generalizable to population of interest
- ▶ More convenient and less costly

Study design

Experimental studies (e.g., RCTs)

- ▶ Researchers directly control exposures or treatments
- ▶ Ability to make causal statements
- ▶ Less real-world applicability and generalizability

Observational studies (e.g., surveys)

- ▶ Researchers do not assign exposures or treatments
- ▶ Real-world setting with lower burden on participants
- ▶ Inability to prove causality

What can go wrong?

Selection bias, reporting bias, non-response bias, attrition bias, spin bias, confounding, detection bias, lack of blinding, straight up falsified data (this happens), ...

Catalogue of Bias



...and so much more.

In recent news...

COVID-19 Antibody Seroprevalence in Santa Clara County, California

Eran Bendavid, Bianca Mulaney, Neeraj Sood, Soleil Shah, Emilia Ling, Rebecca Bromley-Dulfano, Cara Lai, Zoe Weissberg, Rodrigo Saavedra-Walker, James Tedrow, Dona Tversky, Andrew Bogan, Thomas Kupiec, Daniel Eichner, Ribhav Gupta, John Ioannidis, Jay Bhattacharya


doi: <https://doi.org/10.1101/2020.04.14.20062463>

This article is a preprint and has not been peer-reviewed [what does this mean?]. It reports new medical research that has yet to be evaluated and so should *not* be used to guide clinical practice.

(we'll be revisiting this study later this semester...)

In recent news...

Peer Review of “COVID-19 Antibody Seroprevalence in Santa Clara County, California”

 Balaji S. Srinivasan [Follow](#)

Apr 17 · 14 min read ★

[Twitter](#) [LinkedIn](#) [Facebook](#) [Bookmark](#)

The high reported positive rate in this serosurvey may be explained by the false positive rate of the test and/or by sample recruitment issues.

What do some other people have to say about this study...?

Reproducibility and replicability

Reproducibility: being able to take the original data and code to reproduce all numerical findings

Replicability: being able to independently repeat an entire study without use of the original data (generally with the same methods)

Some best practices from the American Statistical Association:

- ▶ End-to-end scripting of research
- ▶ Use of version control and documentation
- ▶ Publication of code along with data

The current replication crisis

Drip, drip: Former Harvard stem cell researcher up to 18 retractions

Piero Anversa, a former star researcher at Harvard Medical School who left the institution under a cloud, is up to 18 retractions. But that's barely half of the 31 papers by Anversa's group that Harvard has requested journals pull over concerns about the integrity of the findings.

The two articles, published in the *Proceedings of the National Academy of Sciences*, appeared in 2008 and 2009. Anversa author, Annarosa Leri, are among the authors

Withdrawal: Maturation of lipoprotein in the endoplasmic reticulum: Control of formation of functional dimers and aggregates.

Osnat Ben-Zeev, Hui Z. Mao and Mark H. Doolittle

VOLUME 277 (2002) PAGES 10727-10738

This article has been withdrawn by Osnat Ben-Zeev and Mark H. Doolittle. Fig. 3A contained several duplicated regions. Figs. 6, B and C; 7, A and B; and 8, A and B, were inappropriately marked.

© 2019 by The American Society for Biochemistry and Molecular Biology



THE UNITED STATES
DEPARTMENT OF JUSTICE

Home » Office of Public Affairs » News

JUSTICE NEWS

Department of Justice

Office of Public Affairs

FOR IMMEDIATE RELEASE

Monday, March 25, 2019

Duke University Agrees to Pay U.S. \$112.5 Million to Settle False Claims Act Allegations Related to Scientific Research Misconduct

RETRACTION | AUGUST 16, 2019

Kohrt HE, Houot R, Goldstein MJ, Weiskopf K, Alizadeh AA, Brody J, Müller A, Pachynski R, Czerwinski D, Coutre S, Chao MP, Chen L, Tedder TF, Levy R. CD137 stimulation enhances the antilymphoma activity of anti-CD20 antibodies. *Blood*. 2011;117(8):2423-2432. [a](#)

Blood (2019) 134 (7): 686.

<https://doi.org/10.1182/blood.201902416>

[Article history](#)

Connected Content

This is a retraction to: [CD137 stimulation enhances the antilymphoma activity of anti-CD20 antibodies](#)

The Editors of *Blood* retract the 24 February 2011 paper cited above. Concerns regarding the data underlying Figures 3A, 3B, 3C, 4A, and 5C were brought to the attention of Stanford University. The university investigated the issue, conducting a search for any original sources of these data. The search was unsuccessful; the experiments, data, and figure preparation for these figures were overseen by Holbrook E. Kohrt, who died before the university became aware of the concerns. As a result, the data underlying these figures cannot be validated.

Course syllabus

The course syllabus is the official document regarding all policies and guidelines and serves as the course syllabus. It is available on the course website [here](#); a .pdf version is available on Sakai.