

- » Introduction
- » Language:
'Prediction' is retrospective
- » Definitions:
'Prediction' is correlation
- » Validity:
Correlations can overfit
- » Paradox:
'Truth' may not predict
- » Summary
- » References

» What Everybody Needs to Know About 'Prediction' in Machine Learning

» *Momin M. Malik, PhD <momin_malik@cyber.harvard.edu>*

Data Science Postdoctoral Fellow

Berkman Klein Center for Internet & Society at Harvard University

Leverhulme Centre for the Future of Intelligence, University of Cambridge
3 December 2018

Slides: <https://mominmalik.com/cfi.pdf>

► Existential threats, or myths?

► Introduction

► Language:
'Prediction' is retrospective

► Definitions:
'Prediction' is correlation

► Validity:
Correlations can overfit

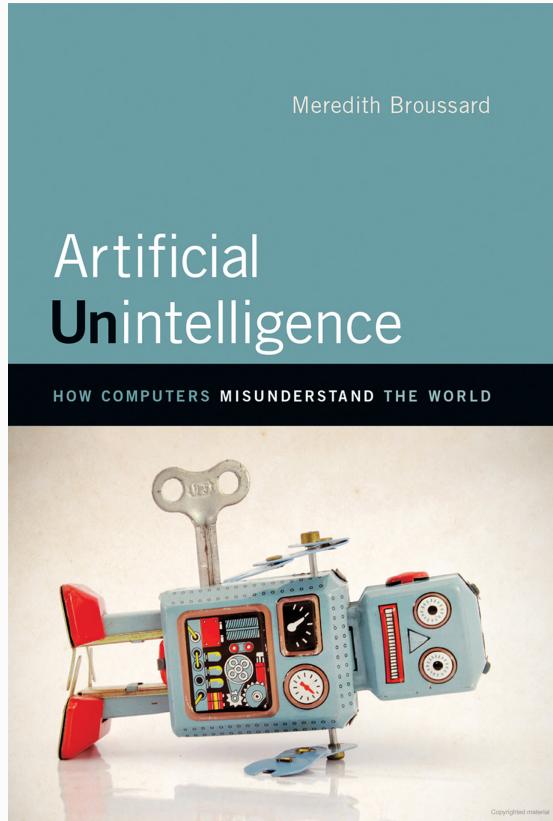
► Paradox:
'Truth' may not predict

► Summary

► References



› Solid general resource



- › Chapter 7: "Machine Learning: The DL on ML"
 - (Two mistakes, see https://github.com/momin-malik/guides/blob/master/Broussard_DL_on_ML.ipynb)
- › Chapter 3: "Hello, AI"
 - "So, it's not real AI?" he asked.
 - "Oh, it's real," I said. "And it's spectacular. But you know, don't you, that there's no simulated person inside the machine? Nothing like that exists. It's computationally impossible."
 - His face fell. "I thought that's what AI meant," he said. "I heard about IBM Watson, and the computer that beat the champion at Go, and self-driving cars. I thought they invented real AI."
- › Then, the rest of Part II

➤ The things everybody needs to know

- Introduction
- Language:
‘Prediction’ is retrospective
- Definitions:
‘Prediction’ is correlation
- Validity:
Correlations can overfit
- Paradox:
‘Truth’ may not predict
- Summary
- References

- Language: ‘Prediction’ (technical term) is not prediction (colloquial term); prediction is prospective, ‘prediction’ is retrospective.
- Definitions: ‘Prediction’ is based on correlations
- Validity: Correlations can *overfit*, and cross-validation only partially addresses
- Paradox: The *bias-variance tradeoff* (a consequence of the definition) makes it possible for a ‘false’ model to predict better than a ‘true’ one



» Introduction

» Language:
'Prediction' is
retrospective

» Definitions:
'Prediction' is
correlation

» Validity:
Correlations
can overfit

» Paradox:
'Truth' may
not predict

» Summary

» References

► Language: 'Prediction' is not prediction

➤ Lots of “predict...”

➤ Introduction

➤ Language:
‘Prediction’ is retrospective

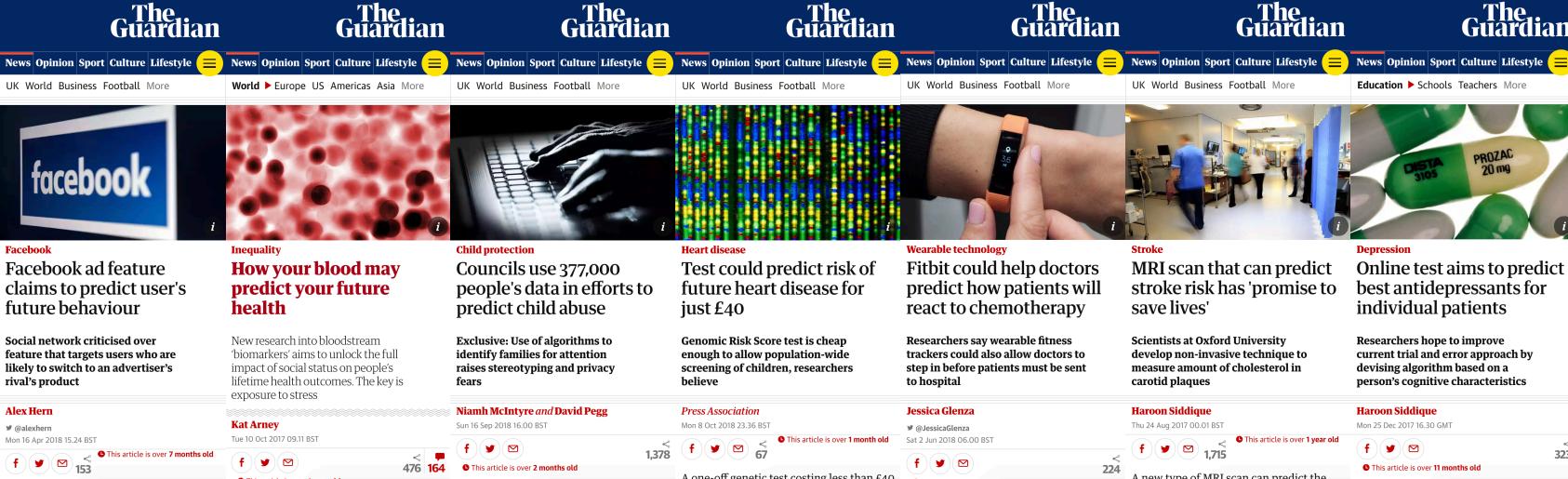
➤ Definitions:
‘Prediction’ is correlation

➤ Validity:
Correlations can overfit

➤ Paradox:
‘Truth’ may not predict

➤ Summary

➤ References



The Guardian headlines include:

- Facebook**: Facebook ad feature claims to predict user's future behaviour
- Inequality**: How your blood may predict your future health
- Child protection**: Councils use 377,000 people's data in efforts to predict child abuse
- Heart disease**: Test could predict risk of future heart disease for just £40
- Wearable technology**: Fitbit could help doctors predict how patients will react to chemotherapy
- Stroke**: MRI scan that can predict stroke risk has 'promise to save lives'
- Depression**: Online test aims to predict best antidepressants for individual patients

Each article includes a small image, a brief summary, and social media sharing options.

➤ If you relied on *The Guardian*, what sort of picture might you get?

➤ 'Prediction' is not prediction!

➤ Introduction

➤ Language:
'Prediction' is retrospective

➤ Definitions:
'Prediction' is correlation

➤ Validity:
Correlations can overfit

➤ Paradox:
'Truth' may not predict

➤ Summary

➤ References

*"I Wanted to Predict Elections with Twitter
and all I got was this Lousy Paper"*

A Balanced Survey on Election Prediction using Twitter Data

Daniel Gayo-Avello
dani@uniovi.es
@PFCdgayo

Department of Computer Science - University of Oviedo (Spain)

May 1, 2012

Abstract

Predicting X from Twitter is a popular fad within the Twitter research subculture. It seems both appealing and relatively easy. Among such kind of studies, electoral prediction is maybe the most attractive, and at this moment there is a growing body of literature on such a topic.

This is not only an interesting research problem but, above all, it is extremely difficult. However, most of the authors seem to be more interested in claiming positive results than in providing sound and reproducible methods.

"It's not prediction at all! I have not found a single paper predicting a future result. All of them claim that a prediction could have been made; i.e. they are post-hoc analysis and, needless to say, negative results are rare to find."

➤ “Wishful mnemonics” of AI

➤ Introduction

➤ Language:
‘Prediction’ is retrospective

➤ Definitions:
‘Prediction’ is correlation

➤ Validity:
Correlations can overfit

➤ Paradox:
‘Truth’ may not predict

➤ Summary

➤ References

ARTIFICIAL INTELLIGENCE MEETS NATURAL STUPIDITY

Drew McDermott

MIT AI Lab Cambridge, Mass 02139

As a field, artificial intelligence has always been on the border of respectability, and therefore on the border of crackpottery. Many critics <Dreyfus, 1972>, <Lighthill, 1973> have urged that we are over the border. We have been very defensive toward this charge, drawing ourselves up with dignity when it is made and folding the cloak of Science about us. On the other hand, in private, we have been justifiably proud of our ideas, because pursuing them is the only

Unfortunately, the necessity for s the culture of the hacker in computer to cripple our self-discipline. In a young field, self-discipline is not necessarily a virtue, but we are not getting any younger. In the past few years, our tolerance of sloppy thinking has led us to repeat many mistakes over and over. If we are to retain any credibility, this should stop.

This paper is an effort to ridicule some of these mistakes. Almost everyone I know should find himself the target at some point or other; if you don't, you are encouraged to write up your own favorite fault. The three described here I suffer from myself. I hope self-ridicule will be a complete catharsis, but I doubt it. Bad

though, if we can't

Wishful Mnemonics

Wishful Mnemonics

A major source of simple-mindedness in AI programs is the use of mnemonics like "UNDERSTAND" or "GOAL" to refer to programs and data structures. This practice has been inherited from more

Compare the mnemonics in Planner <Hewitt,1972> with those in Conniver <Sussman and McDermott, 1972>:

Planner	Conniver
GOAL	FETCH & TRY-NEXT
CONSEQUENT	IF-NEEDED
ANTECEDENT	IF-ADDED
THEOREM	METHOD
ASSERT	ADD

It is so much harder to write programs using the terms on the right! When you say (GOAL . . .), you can just feel the enormous power at your fingertips. It is, of course, an illusion.

When you say (GOAL . . .), you can just feel the enormous power at your fingertips. It is, of course, an illusion.

1965> What if atomic symbols had been called "concepts", or CONS had been called ASSOCIATE? As it is, the programmer has no debts to pay to the system. He can build whatever he likes. There are some minor faults; "property lists" are a little risky; but by now the term is sanitized.

Resolution theorists have been pretty good about wishful mnemonics. They thrive on hitherto meaningless words like RESOLVE and PARAMODULATE, which can only have their humble, technical meaning. There are actually quite few pretensions in the resolution literature. <Robinson, 1965> Unfortunately, at the top of their intellectual edifice stand the word "deduction". This is very wishful, but not entirely their fault. The logicians who first misused the term (e.g., in the "deduction" theorem) didn't have our problems; pure resolution theorists don't either. Unfortunately, too many AI researchers took them at their word and assumed that deduction, like payroll processing, had been tamed.

Of course, as in many such cases, the only consequence in the long run was that "deduction" changed in meaning, to become something narrow, technical, and not a little sordid.

➤ Proposal: More precise language

- Introduction
- Language:
'Prediction' is retrospective
- Definitions:
'Prediction' is correlation
- Validity:
Correlations can overfit
- Paradox:
'Truth' may not predict
- Summary
- References

- ~~Predict the risk, predict the likelihood~~: Calculate the risk, calculate the likelihood
- ~~Predict the probability~~: Estimate the probability
- ~~Prediction, predicted~~: Fitted value, fitted
- ~~We predict~~: We detect, we classify, we model
- ~~X predicts Y~~: X is correlated with Y
- ~~X predicts Y, ceteris paribus~~ (partial correlation): X is associated with Y

➤ Proposal: Use alternatives

- Introduction
- Language:
'Prediction' is retrospective
- Definitions:
'Prediction' is correlation
- Validity:
Correlations can overfit
- Paradox:
'Truth' may not predict
- Summary
- References
- Retrodiction
- Backtesting (retrodiction for testing)
- Hindcasting (backtesting for forecasting)
- In-sample vs. ➤ Out of-sample
- Interpolation vs. ➤ Extrapolation
- Diagnosis vs. ➤ Prognosis
- Retrospective vs. ➤ Prospective

➤ (Language not enough: *mechanics matter*)

➤ Introduction

➤ Language:
'Prediction' is retrospective

➤ Definitions:
'Prediction' is correlation

➤ Validity:
Correlations can overfit

➤ Paradox:
'Truth' may not predict

➤ Summary

➤ References

Pseudo-Mathematics and Financial Charlatanism: The Effects of Backtest Overfitting on Out-of-Sample Performance

*David H. Bailey, Jonathan M. Borwein,
Marcos López de Prado, and Qiji Jim Zhu*

Another thing I must point out is that you cannot prove a vague theory wrong. [...] Also, if the process of computing the consequences is indefinite, then with a little skill any experimental result can be made to look like the expected consequences

"training set" in the machine-learning literature). The OOS performance is simulated over a sample not used in the design of the strategy (a.k.a. "testing set"). A backtest is *realistic* when the IS performance

› Introduction

› Language:
'Prediction' is retrospective› Definitions:
'Prediction' is correlation› Validity:
Correlations can overfit› Paradox:
'Truth' may not predict

› Summary

› References

► Definitions: 'Prediction' is correlation, not causation

► Prediction is correlation

- Introduction
- Language: 'Prediction' is retrospective
- Definitions: 'Prediction' is correlation
- Validity: Correlations can overfit
- Paradox: 'Truth' may not predict
- Summary
- References

- > Prediction = "Fitted value" minimizing *loss*
- > $L(y, f(x)) = (y - f(x))^2$
- > Spurious (non-causal) correlations can *fit* really well!
- > But such fits fall apart if the context changes (Google Flu Trends)

POLICYFORUM

BIG DATA

The Parable of Google Flu: Traps in Big Data Analysis

David Lazer,^{1,2*} Ryan Kennedy,^{1,3*} Gary King,² Alessandro Vespiagnani^{1,6,3}

In February 2013, Google Flu Trends (GFT) made headlines but not for a reason that Google executives or the creators of the flu tracking system would have hoped. *Nature* reported that GFT was predicting more than double the proportion of doctor visits for influenza-like illness (ILI) than the Centers for Disease Control and Prevention (CDC), which bases its estimates on surveillance reports from laboratories across the United States (1, 2). This happened despite the fact that GFT was built to predict CDC reports. Given that GFT is often held up as an exemplary use of big data (3, 4), what lessons can we draw from this error?

The problems we identify are not limited to GFT. Research on whether search or social media can predict x has become commonplace (5–7) and is often put in sharp contrast with traditional methods and hypotheses. Although these studies have shown the value of these data, we are far from a place where they can supplant more traditional methods or theories (8). We explore two issues that contributed to GFT's mistakes—big data hubris and algorithm dynamics—and offer lessons for moving forward in the big data age.

Big Data Hubris

"Big data hubris" is the often implicit assumption that big data are a substitute for, rather than a supplement to, traditional data collection and analysis. Elsewhere we



ability and dependencies among data (12). The core challenge is that most big data that have received popular attention are not the output of instruments designed to produce valid and reliable data amenable for scientific analysis.

The initial version of GFT was a particularly problematic marriage of big and small data. Essentially, the methodology was to find the best matches among 50 million search terms to fit 1152 data points (13). The odds of finding search terms that match the propensity of the flu but are structurally unrelated, and so do not predict the future, were quite high. GFT developers, in fact, report *wading* through seasonal search

Large errors in flu prediction were largely avoidable, which offers lessons for the use of big data.

run ever since, with a few changes announced in October 2013 (10, 15).

Although not widely reported until 2013, the new GFT has been persistently overestimating flu prevalence for a much longer time. GFT also missed by a very large margin in the 2011–2012 flu season and has missed high for 100 out of 108 weeks starting with August 2011 (see the graph). These errors are not randomly distributed. For example, last week's errors predict this week's errors (temporal autocorrelation), and the direction and magnitude of error varies with the time of year (seasonality). These patterns mean that GFT overlooks considerable information that could be extracted by traditional statistical methods.

Even after GFT was updated in 2009, the comparative value of the algorithm as a stand-alone flu monitor is questionable. A study in 2010 demonstrated that GFT accuracy was not much better than a fairly simple projection forward using already available (typically on a 2-week lag) CDC data (4). The comparison has become even worse since that time, with lagged models significantly outperforming GFT (see the graph). Even 3-week-old CDC data do a better job of projecting current flu prevalence than GFT [see supplementary materials (SM)].

Considering the large number of approaches that provide inference on influenza activity (16–18), does this mean that

> (Caution: “Causation” is itself limited)

- > Introduction
- > Language: ‘Prediction’ is retrospective
- > Definitions: ‘Prediction’ is correlation
- > Validity: Correlations can overfit
- > Paradox: ‘Truth’ may not predict
- > Summary
- > References

- > Critique 1: Causal inference (econometrics) can fail hopelessly
- > Critique 2: Automated methods (from “causal learning”) have strong, unrealistic, and untestable assumptions
- > Critique 3: Statistical expression of causation is short-range (Gene Richardson)

Sociological Methods & Research
39(2) 258–282
© The Author(s) 2010
Reprints and permission:
[sagepub.com/journalsPermissions.nav](http://sagepub.com)
DOI: 10.1177/0049124110378998
<http://smr.sagepub.com>
\$SAGE

A Cautionary Note on the Use of Matching to Estimate Causal Effects: An Empirical Example Comparing Matching Estimates to an Experimental Benchmark

Kevin Arceneaux¹, Alan S. Gerber², and Donald P. Green²

Abstract

In recent years, social scientists have increasingly turned to matching as a method for drawing causal inferences from observational data. Matching compares those who receive a treatment to those with similar background attributes who do not receive a treatment. Researchers who use matching frequently tout its ability to reduce bias, particularly when applied to data sets that contain extensive background information. Drawing on a randomized voter mobilization experiment, the authors compare estimates generated by matching to an experimental benchmark. The enormous sample size enables the authors to exactly match each treated subject to 40 untreated subjects. Matching greatly exaggerates the effectiveness of pre-election phone calls encouraging voter participation. Moreover, it can produce nonsensical results: Matching suggests that another pre-election phone

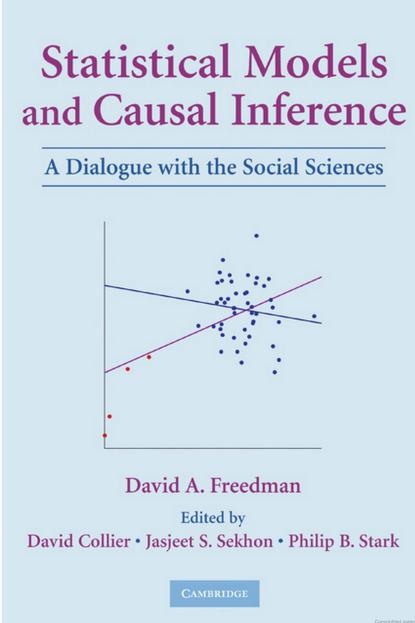
¹Temple University, Philadelphia, PA, USA

²Yale University, New Haven, CT, USA

Corresponding Author:

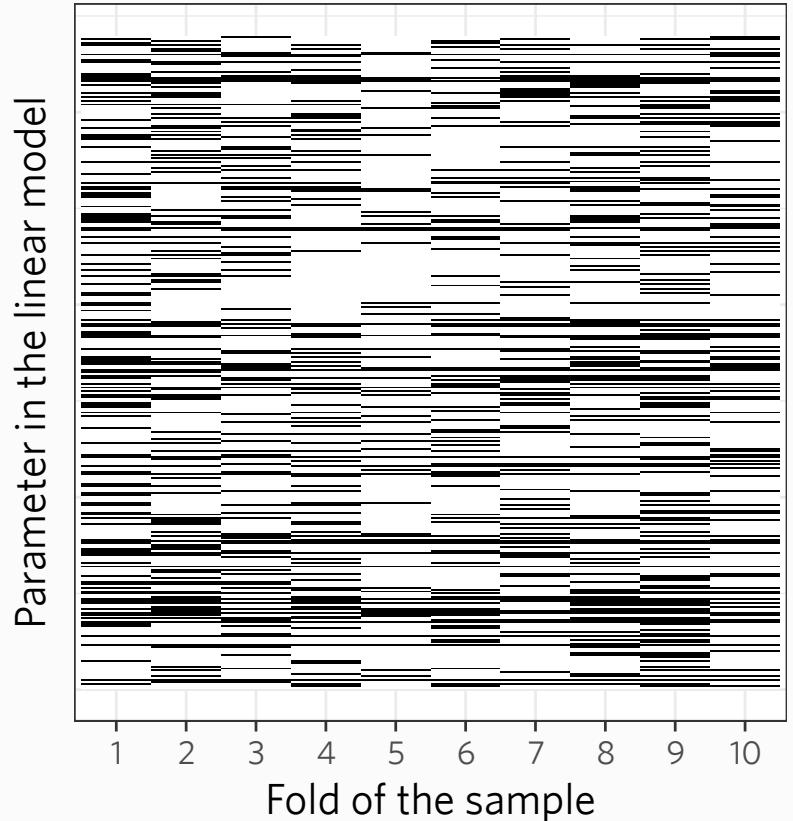
Kevin Arceneaux, Department of Political Science and Institute for Public Affairs, Temple University, 453 Gladfelter Hall, 1115 West Berks St., Philadelphia, PA 19122, USA
Email: kevin.arceneaux@temple.edu

Downloaded from smr.sagepub.com at CARNEGIE MELLON UNIV LIBRARY on October 13, 2015



➤ The problem with correlation

- Very different models will 'predict' equally well, and often better than any theory-driven model (Mullainathan & Spiess, 2017)
- For *intervention*, we need causality (or at least associations)
- Another problem: correlations can *overfit*



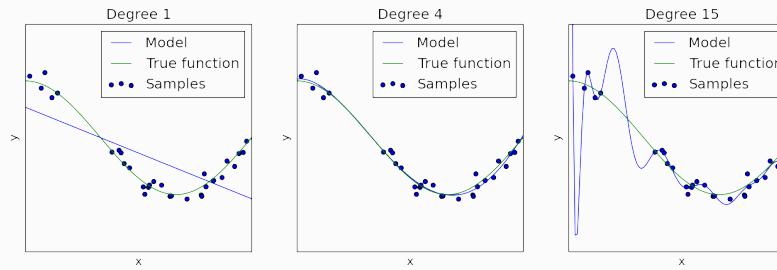
- Introduction
- Language:
'Prediction' is retrospective
- Definitions:
'Prediction' is correlation
- Validity:
Correlations can overfit
- Paradox:
'Truth' may not predict
- Summary
- References

➤ **Validity: Correlations can overfit, cross-validation doesn't fully address**

Overfitting and cross validation

- ▶ Introduction
- ▶ Language: 'Prediction' is retrospective
- ▶ Definitions: 'Prediction' is correlation
- ▶ Validity: Correlations can overfit
- ▶ Paradox: 'Truth' may not predict
- ▶ Summary
- ▶ References

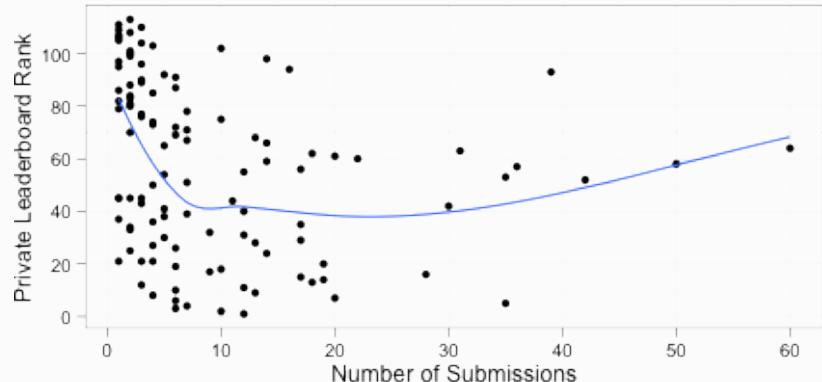
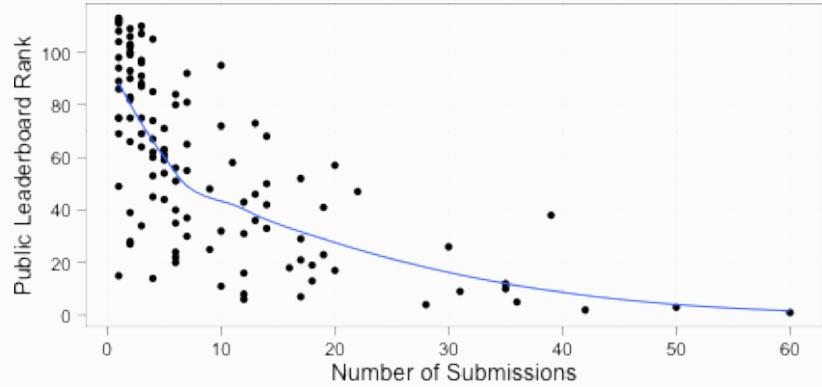
- ▶ Overfitting: Model correlates with the 'noise' rather than causes/*data-generating process*. Machine learning metaphor: "memorize the data."



- ▶ (A *p*-value compares fit and *uncertainty*; overfitting is simpler)
- ▶ Cross validation: split the data into two parts (e.g., 1:1, 4:1, 9:1). *The central tendency should be the same, but not the noise.* Error rate on the held-out "test" set should say how well correlations generalize

➤ But cross-validation can fail

- Re-using the test set can overfit to the test set!
Happens in Kaggle
- Or, if there are dependencies (temporal, network, group) between data splits, it “shares” information
- E.g., temporal: Fitting on values that come after test values is “time traveling”!



› Introduction

› Language:
'Prediction' is
retrospective› Definitions:
'Prediction' is
correlation› Validity:
Correlations
can overfit› Paradox:
'Truth' may
not predict

› Summary

› References

➤ A 'false' model may predict better than a 'true' one

➤ The bias-variance tradeoff

- The bias-variance ‘decomposition’, a foundational result for machine learning and modern statistics:

$$\begin{aligned}\text{EPE}(x) &= \mathbb{E}[(Y - \hat{f}(x))^2 | X = x] \\ &= \text{Var}(Y) + \mathbb{E}[(\hat{f}(x) - f(x))^2 | X = x] + \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2 | X = x] \\ &= \sigma^2 + \text{bias}^2(\hat{f}(x)) + \text{Var}(\hat{f}(x))\end{aligned}$$

- Leads to a ‘tradeoff’: *Even if we have all the “right” variables, a biased model may be better*
- This is very strange!

► Simulation illustration: Setup

- A linear data-generating process.

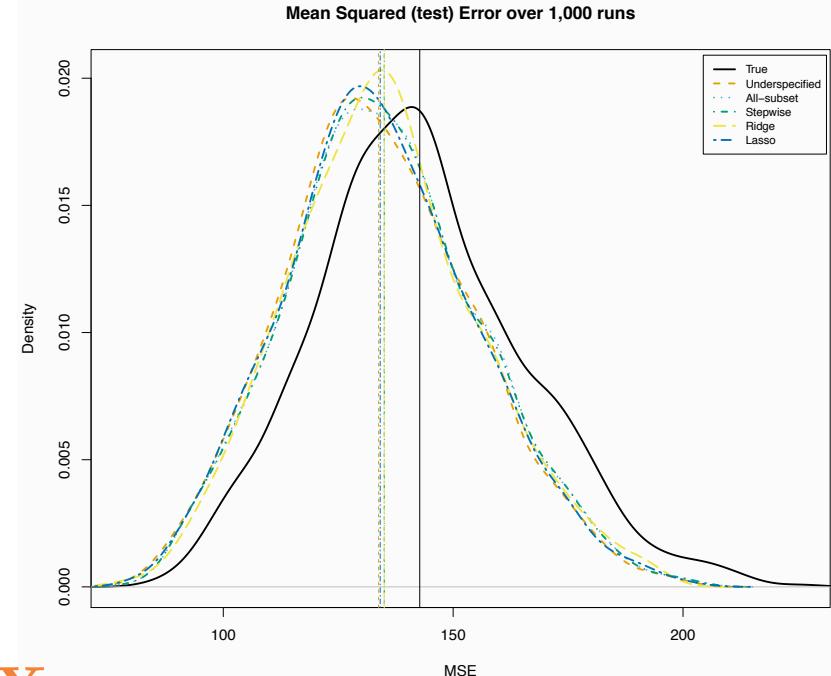
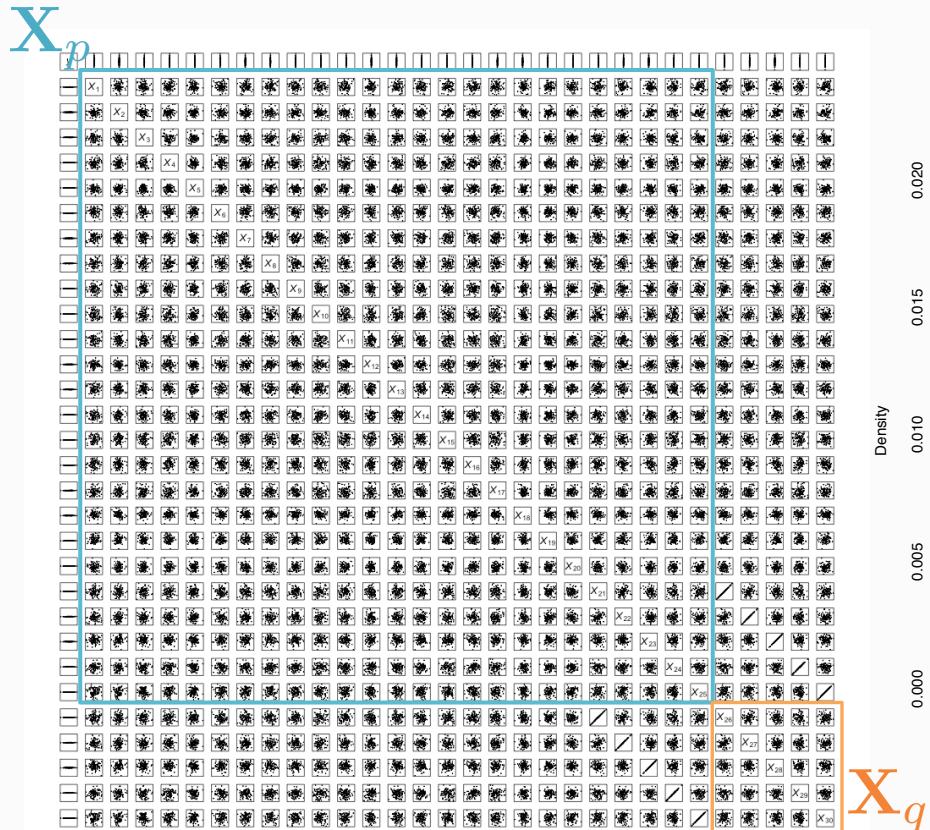
$$\mathbf{y} \sim \mathcal{N}(\beta_p \mathbf{X}_p + \beta_q \mathbf{X}_q, \sigma^2 \mathbf{I})$$

- Wu et al. (2007): Fitting only \mathbf{X}_p has lower expected MSE than fitting the model that generated the data when:

$$\beta_q^T \mathbf{X}_q^T (\mathbf{I}_n - \mathbf{H}_p) \mathbf{X}_q \beta_q < q\sigma^2$$

► The 'true' model predicts worse!

- Introduction
- Language:
'Prediction' is retrospective
- Definitions:
'Prediction' is correlation
- Validity:
Correlations can overfit
- Paradox:
'Truth' may not predict
- Summary
- References



➤ Summary

- Introduction
- Language:
‘Prediction’ is retrospective
- Definitions:
‘Prediction’ is correlation
- Validity:
Correlations can overfit
- Paradox:
‘Truth’ may not predict
- Summary
- References

- ‘Prediction’ is a metaphor used for fitted values, not (necessarily) actual prediction
- Spurious correlations count as ‘prediction’ and can do quite well in narrow terms, but are fragile and don’t help us intervene
- Correlations can overfit, and cross-validation doesn’t fully solve
- The bias-variance tradeoff means things are even more strange
- *I would argue: These are the pertinent issues*

› Thank you!

› Introduction

› Language:
'Prediction' is
retrospective

› Definitions:
'Prediction' is
correlation

› Validity:
Correlations
can overfit

› Paradox:
'Truth' may
not predict

› Summary

› References



► Citations/Credits by slide number

► Introduction

► Language:
'Prediction' is retrospective

► Definitions:
'Prediction' is correlation

► Validity:
Correlations can overfit

► Paradox:
'Truth' may not predict

► Summary

► References

- 2 Robot holding skull: Cover image of "What Will Become of Us?", *New York Times Magazine* (The Tech & Design issue), 14 November 2018. Concept by delcan & company. Photo illustration by Jamie Chung. Prop styling by Pink Sparrow. C.G. work by Justin Metz.
<https://www.nytimes.com/2018/11/14/magazine/behind-the-cover-what-will-become-of-us.html>.
- 2 Terminator skull: Nemesis Now Ltd, Terminator Skull Box T-800 (18CM).
<https://www.menkind.co.uk/terminator-t800-skull-box>.
- 2 Hand: Hayati Kayhan, Holding human skull in hand, Conceptual image (Shakespeare's Hamlet scene concept). 19 October 2014.
- 3 Meredith Broussard, *Artificial Unintelligence: How Computers Misunderstand the World*. MIT Press, 2018.
- 7 Sitaram Asur and Bernardo A. Huberman, "Predicting the Future With Social Media." In *Proceedings of the 2010 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology* (WI-IAT '10), 492–499. 2010. <https://dx.doi.org/10.1109/WI-IAT.2010.63>.
- 7 Ziad Obermeyer and Ezekiel J. Emanuel, "Predicting the Future — Big Data, Machine Learning, and Clinical Medicine." *New England Journal of Medicine* 375, no. 13 (2016): 1216–1219. <https://dx.doi.org/10.1056/NEJMp1606181>.
- 7 OED Online, "predict, v." Oxford University Press, July 2018.
<http://www.oed.com/view/Entry/149856>.
- 8 Daniel Gayo-Avello, "'I Wanted to Predict Elections with Twitter and all I got was this Lousy Paper': A Balanced Survey on Election Prediction using Twitter Data." *arXiv*, 28 April 2012. <https://arxiv.org/abs/1204.6441>.
- 8 Daniel Gayo-Avello, "No, You Cannot Predict Elections with Twitter." *IEEE Internet Computing* 16 (2012): 91–94. <https://dx.doi.org/10.1109/MIC.2012.137>.
- 9 Drew McDermott, "Artificial Intelligence meets Natural Stupidity." *SIGART Bulletin* 57 (April 1976): 4–9.
- 12 David H. Bailey, Jonathan M. Borwein, Marcos López de Prado, and Qiji Jim Zhu. "Pseudo-Mathematics and Financial Charlatanism: The Effects of Backtest Overfitting on Out-of-Sample Performance." *Notices of the AMS* 61, no. 5 (2014): 458–471. <https://dx.doi.org/10.1090/noti1105>.
- 14 David Lazer, Ryan Kennedy, Gary King, and Alessandro Vespignani, "The Parable of Google Flu Trends: Traps in Big Data Analysis." *Science* 343 (14 March 2014): 1203–1205. <https://dx.doi.org/10.1126/science.1248506>.
- 15 Leo Breiman, Statistical Modeling: The Two Cultures (with Comments and a Rejoinder by the Author). *Statistical Science* 16, no. 3 (2001): 199–231. <https://dx.doi.org/10.1214/ss/1009213726>.
- 15 Galit Shmueli, To Explain or to Predict? *Statistical Science* 25, no. 3 (2010): 289–310. <https://dx.doi.org/10.1214/10-STS330>.
- 15 Sendhil Mullainathan and Jan Spiess, "Machine Learning: An Applied Econometric Approach." *Journal of Economic Perspectives* 31, no. 2 (2017): 87–106. <https://dx.doi.org/10.1257/jep.31.2.87>.
- 16 Kevin Arceneaux, Alan S. Gerber, and Donald P. Green, "A Cautionary Note on the Use of Matching to Estimate Causal Effects: An Empirical Example Comparing Matching Estimates to an Experimental Benchmark." *Sociological Methods & Research* 39, no. 2 (2010): 256–282. <https://dx.doi.org/10.1177/0049124110378098>.
- 16 David A. Freedman. *Statistical Models and Causal Inference: A Dialogue with the Social Sciences*. Cambridge University Press, 2009.
- 19 scikit-learn developers, "Underfitting vs. Overfitting." 2014. https://scikit-learn.org/0.15/auto_examples/plot_underfitting_overfitting.html.
- 20 Greg Park, The Dangers of Overfitting: A Kaggle Postmortem. 6 July 2012. <http://gregpark.io/blog/Kaggle-Psychopathy-Postmortem>.
- 23 Shaohua Wu, T. J. Harris, and K. B. McAuley, "The Use of Simplified or Misspecified Models: Linear Case." *The Canadian Journal of Chemical Engineering* 85, no. 4 (2007): 386–398. <https://dx.doi.org/10.1002/cjce.5450850401>.