

A critical perspective on measurement in digital trace data and machine learning, and implications for demography

Momin M. Malik

Senior Data Science Analyst - AI Ethics, Mayo Clinic

Fellow, Institute in Critical Quantitative, Computational, & Mixed Methodologies

Instructor, University of Pennsylvania School of Social Policy & Practice

Tuesday, 26 April 2022

Max Planck Institute for Demographic Research



MAX-PLANCK-INSTITUT
FÜR DEMOGRAFISCHE
FORSCHUNG

MAX PLANCK INSTITUTE
FOR DEMOGRAPHIC
RESEARCH



Simon Weckert, "Google Maps Hack"

Introduction

Brief historical tour

Bias in geotagged tweets

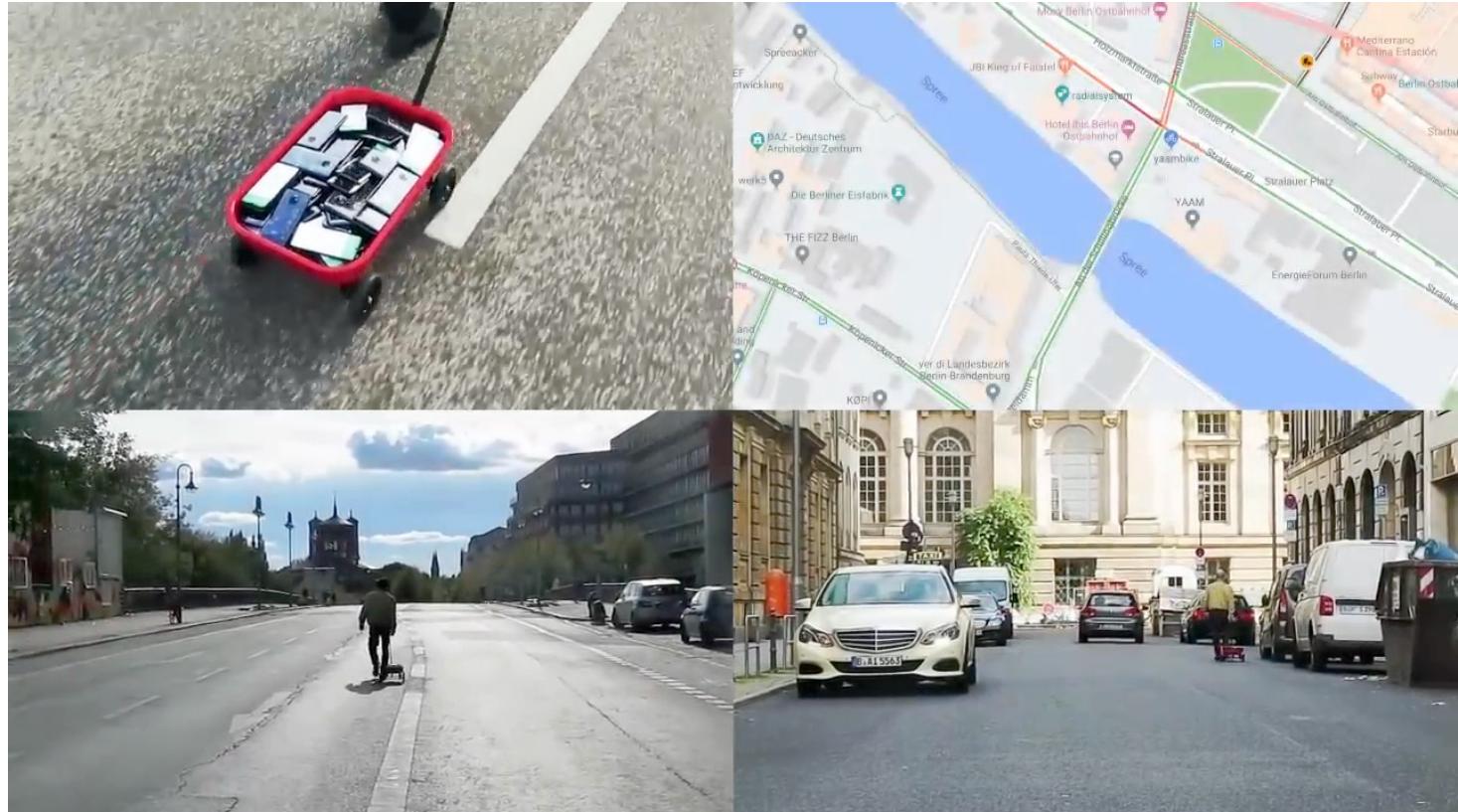
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References





This shows larger themes

- Available data are often only a *proxy*
- So long as the proxy is never the thing itself, it can fail
 - But by interrogating proxies, especially ones we did not construct, we can better understand them
- *Models* of relationships and processes, too, are not the things themselves
- Box (1979): “[For] a model there is no need to ask the question ‘Is the model true?’. If ‘truth’ is to be the ‘whole truth’ the answer must be ‘No’. The only question of interest is ‘Is the model illuminating and useful?’”

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Quick survey

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

- How many people know of Savage and Burrows (2007)? Breiman (2001)?
- What disciplinary backgrounds?
 - Computer science?
 - Statistics? (Math/economics?)
 - Social science?
- How much do you know what machine learning is (or use it)?
 - How is it different from statistics?



Goals and outline

- Brief historical tour: Savage and Burrows (2007) and Breiman (2001)
 - About me
- Bias in geotagged tweets (ICWSM-2015 SPSM)
- Platform effects (ICWSM-2016)
- Hierarchy of limitations in machine learning (2020)
- Problems of cross-validation
- Summary and conclusion

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Introduction

Brief historical tour

Bias in
geotagged
tweets

Platform
effects

Hierarchy of
limitations in
machine
learning

Problems of
cross-
validation

Summary and
conclusion

References

Brief historical tour



Two key historical pieces

- Savage & Burrows (2007): “The coming crisis of empirical sociology”
 - Before Anderson’s “End of theory” (2008) and Lazer et al.’s “Computational social science” (2009)
- Breiman (2001): “Statistical modeling: The two cultures”
 - Even earlier
 - Includes seeds of things we aren’t even fully talked about yet: from problems with interpretability, to limits of cross-validation, to multiplicity of models

Introduction

Brief historical tour

Bias in geotagged tweets

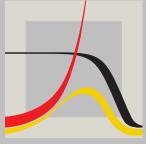
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



"Coming crisis of empirical sociology" (2007)

"In 2004, [Savage] was enrolled in a [ESRC Research Methods festival] session designed to popularize social network methods. He talked about an ESRC-funded research project [on volunteer organizations]... **a postal questionnaire had been sent to 320 members in total**, with a very high response rate. Many members had been interviewed face-to-face to ask detailed questions about their social networks... The resulting intensive study of the members' social ties was amongst the most detailed ever carried out in the UK."

Introduction

Brief historical tour

Bias in geotagged tweets

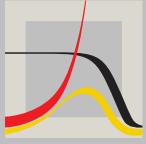
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



"Coming crisis of empirical sociology" (2007)

"During the Festival Savage talked to other participants interested in social network methods. It turned out that one enthusiast was not an academic but worked in a research unit attached to a leading telecommunications company. **When asked what data he used for his social network studies, he shyly replied that he had the entire records of every phone call made on his system over several years, amounting to several billion ties.**"

Introduction

Brief historical tour

Bias in geotagged tweets

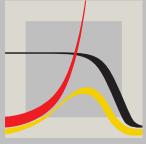
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



"Statistical modeling: The two cultures" (2001)

"the focus in the statistical community on data models has:

- "Led to irrelevant theory and questionable scientific conclusions;
- "Kept statisticians from using more suitable algorithmic models;
- "Prevented statisticians from working on exciting new problems"

"In the past fifteen years, the growth in algorithmic modeling applications and methodology has been rapid. **It has occurred largely outside statistics in a new community—often called machine learning** that is mostly young computer scientists (Section 7). The advances, particularly over the last five years, have been startling."

Introduction

Brief historical tour

Bias in geotagged tweets

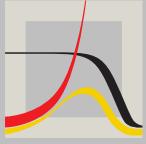
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



“Statistical modeling: The two cultures” (2001)

“Perhaps the damaging consequence of the insistence on data models is that **statisticians have ruled themselves out of some of the most interesting and challenging statistical problems** that have arisen out of the rapidly increasing ability of computers to store and manipulate data. These problems are increasingly present in many fields, both scientific and commercial, and solutions are being found by nonstatisticians.”

“Over the last ten years, there has been a noticeable move toward statistical work on real world problems and reaching out by statisticians toward collaborative work with other disciplines. I believe this trend will continue and, in fact, has to continue **if we are to survive** as an energetic and creative field.”

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



About me

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

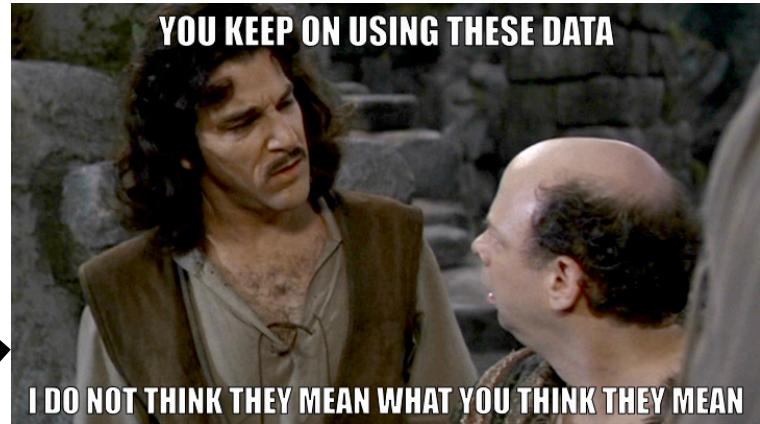
Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

- UG:  DEPARTMENT OF THE HISTORY OF SCIENCE HARVARD UNIVERSITY
- MSc:  
- PhD: Carnegie Mellon University School of Computer Science
 - During:  Data Science For Social Good
Summer Fellowship
- Post-doc:  BERKMAN KLEIN CENTER FOR INTERNET & SOCIETY AT HARVARD UNIVERSITY
- Previously:  AVANT-GARDE HEALTH
- Currently:  Center for Digital Health | Penn Social Policy & Practice UNIVERSITY OF PENNSYLVANIA | ICQCM CRITICAL DATA SCIENCE FOR A DIVERSE WORLD





Introduction

Brief historical tour

**Bias in
geotagged
tweets**

Platform effects

Hierarchy of
limitations in
machine
learning

Problems of
cross-
validation

Summary and
conclusion

References

Bias in geotagged tweets

Momin M. Malik, Hemank Lamba, Constantine Nakos, and Jürgen Pfeffer. 2015. Population bias in geotagged tweets. In *Papers from the 2015 ICWSM Workshop on Standards and Practices in Large-Scale Social Media Research (ICWSM-15 SPSM)*, pages 18–27. May 26, 2015, Oxford, UK. Updated version (2018):
https://www.mominmalik.com/malik_chapter1.pdf



Many maps just show population

Introduction

Brief historical tour

Bias in geotagged tweets

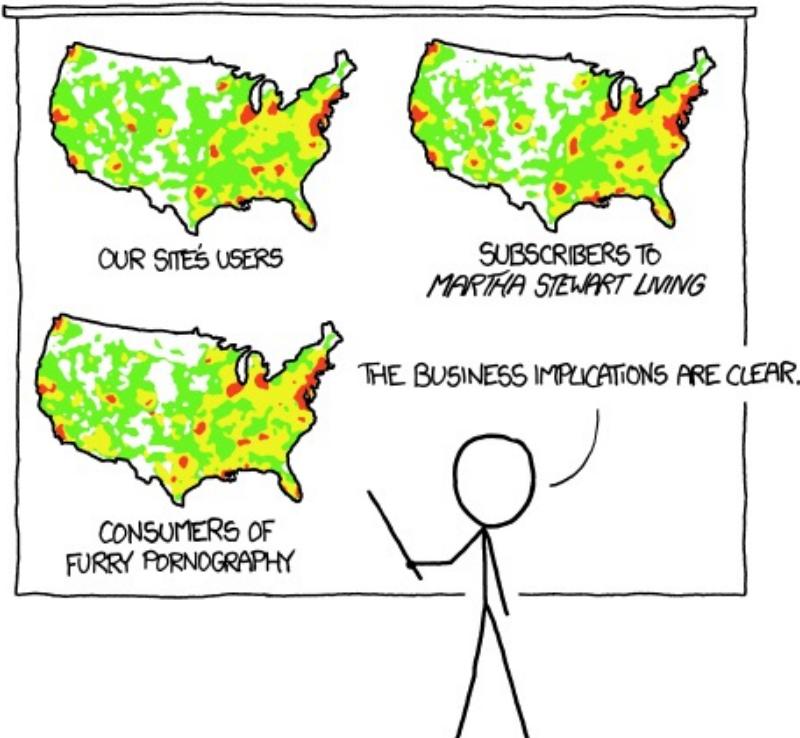
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Randall Munroe. 2012. Heatmap. <https://xkcd.com/1138/>



But maybe we can use this?

Introduction

Brief historical tour

**Bias in
geotagged tweets**

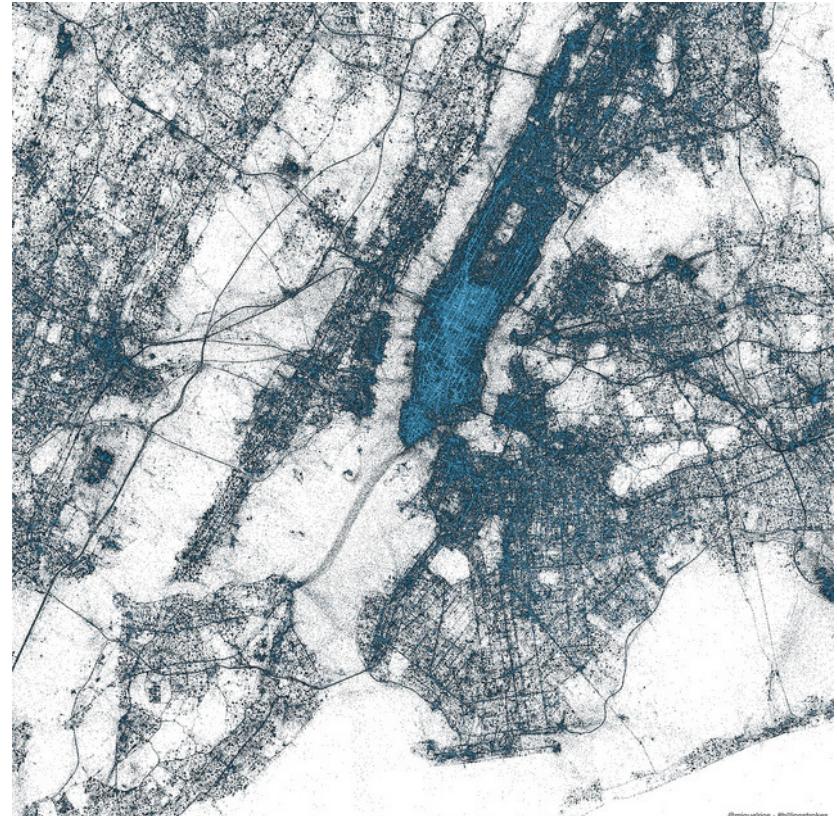
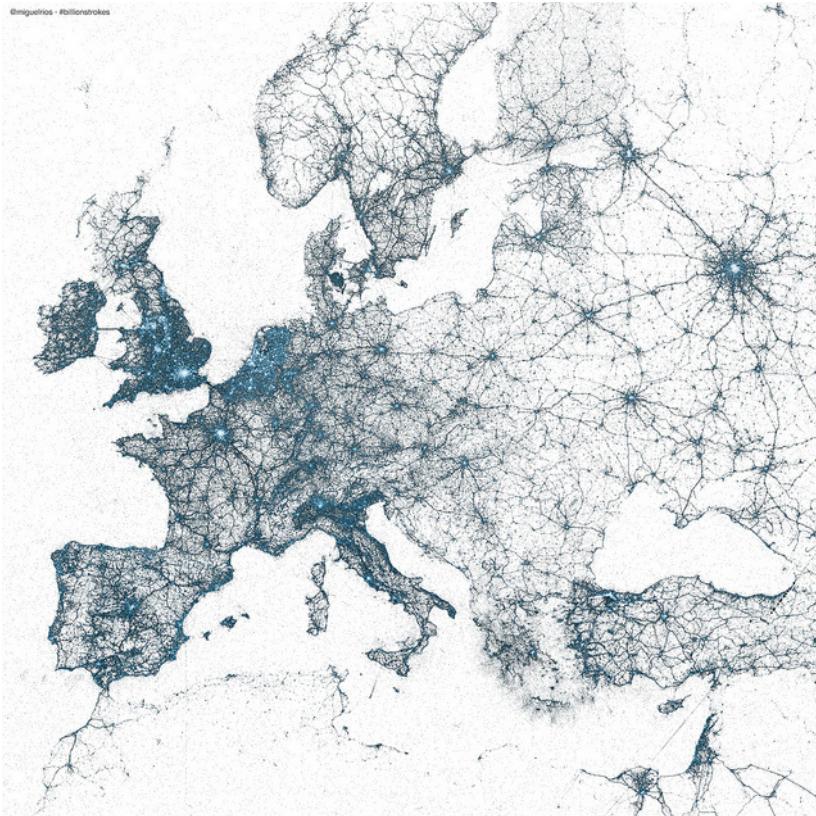
Platform effects

Hierarchy of
limitations in
machine
learning

Problems of
cross-
validation

Summary and
conclusion

References





Do tweets measure population?

Introduction

Brief historical tour

**Bias in
geotagged
tweets**

Platform effects

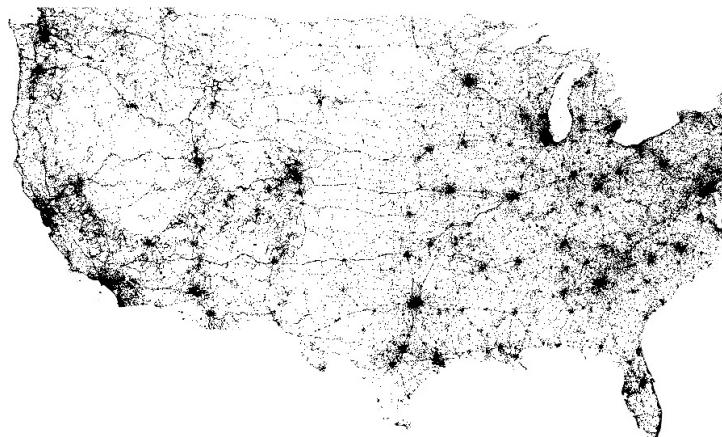
Hierarchy of
limitations in
machine
learning

Problems of
cross-
validation

Summary and
conclusion

References

Geotagged tweets



Adapted from Eric Fischer, 2009, Contiguous United States geotag map.
<https://flic.kr/p/a7WMWS>.

Population



Population density in 2010 US Census. Each square represents 1,000 people. Adapted from Geography Division, U.S. Department of Commerce / Economics and Statistics Administration / U.S. Census Bureau, Nighttime Population Distribution Wall Map.



Modeling population vs. users

- Users proportional to population:

$$U_i = \alpha P_i + \varepsilon_i P_i$$

- Take a log transformation (+Taylor):

$$\log U_i = \log \alpha + \log P_i + \varepsilon'_i$$

- Compare to a linear model:

$$\log U_i = \beta_0 + \beta_1 \log P_i + \varepsilon'_i$$

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

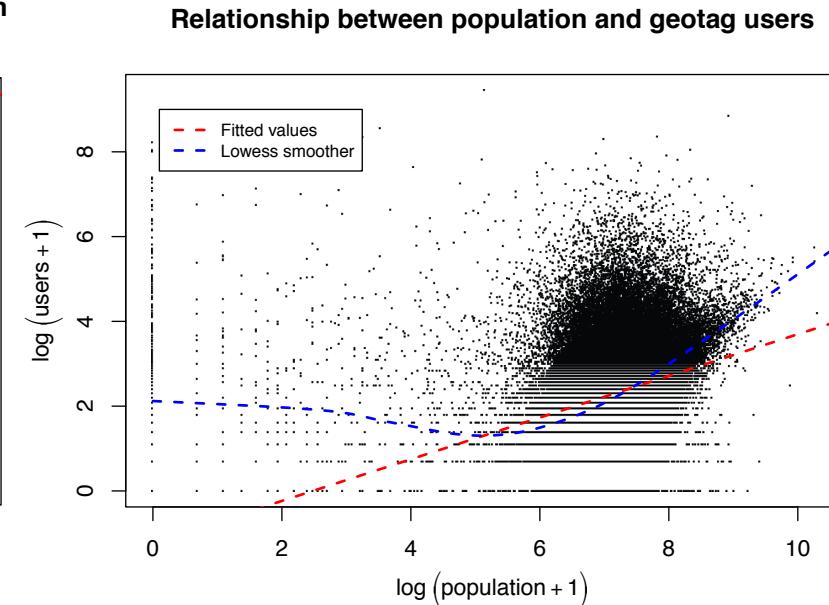
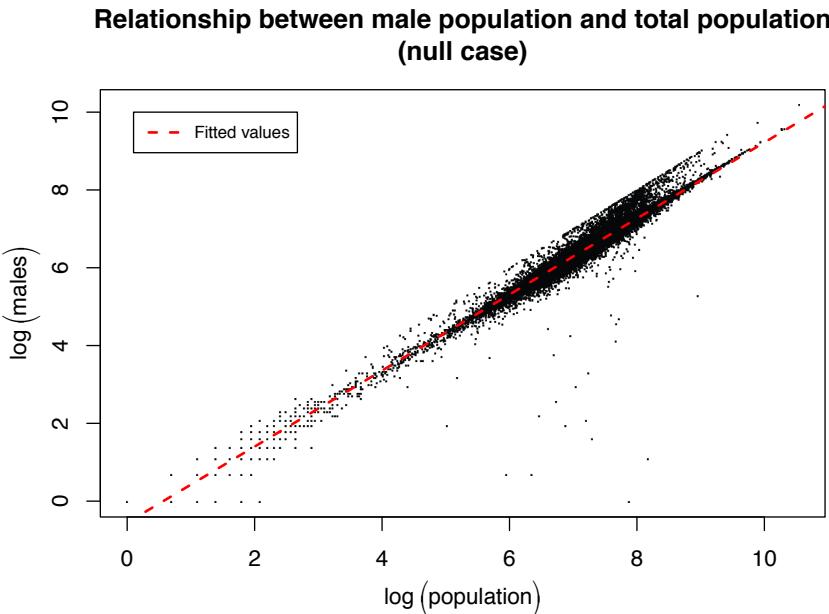
Summary and conclusion

References



Result: Not proportional

(Each dot is a Census block group)



Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Identifying specifics

- Spatial multivariate modeling of biases
Geotagged tweet users associated with:
 - Rural, poor, elderly, non-coastal
 - Asian, Hispanic, black
- ...but these are only the demographics we can access. E.g., harassment of women on Twitter likely discourages geotag use

Introduction

Brief historical tour

**Bias in
geotagged
tweets**

Platform effects

Hierarchy of
limitations in
machine
learning

Problems of
cross-
validation

Summary and
conclusion

References



Why it matters: Some uses are bad

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

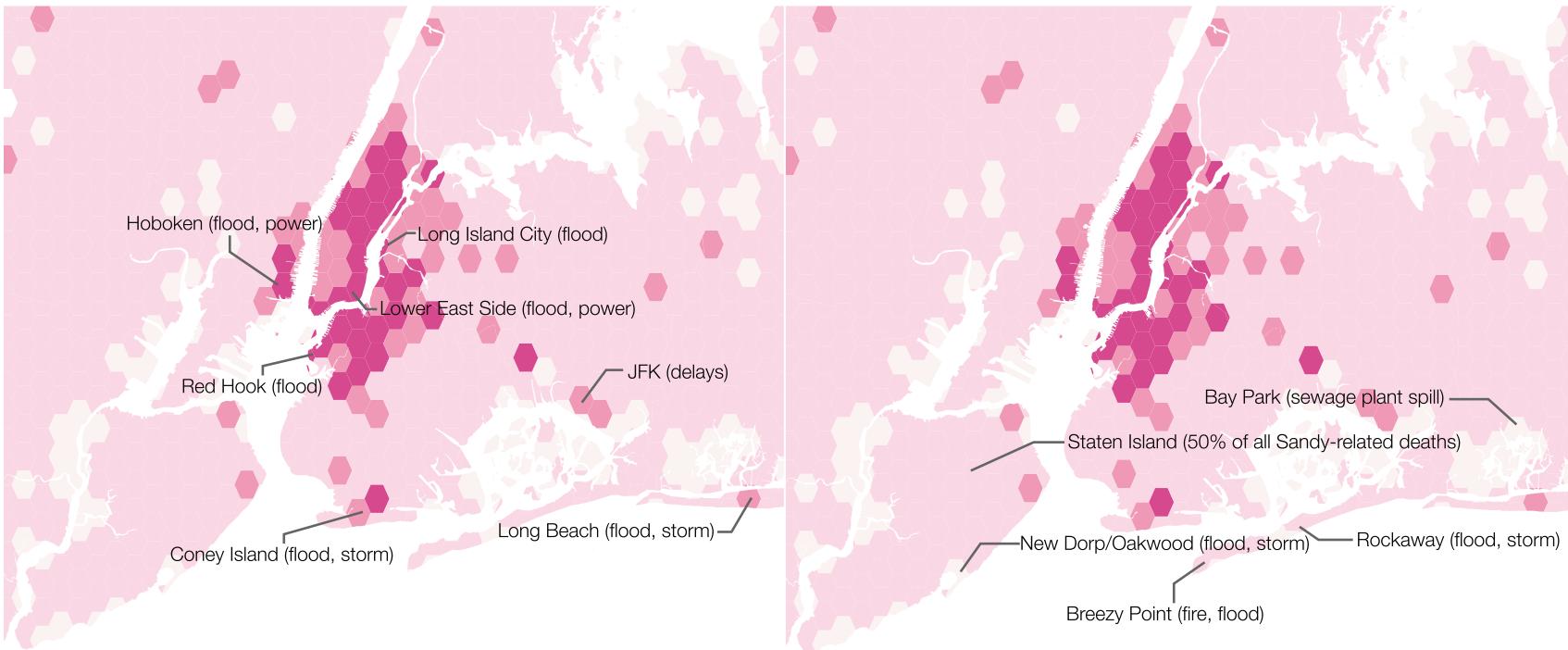
Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

Hurricane Sandy, tweets vs. damage/deaths



Taylor Shelton, Ate Poorthuis, Mark Graham, and Matthew Zook. 2014. Mapping the data shadows of Hurricane Sandy. *Geoforum* 52, 167–179.



Responses to demographic bias

- Model the specific biases!
- Calibration and weighting (Zagheni & Weber, 2015)
- **Use data for appropriate questions**
 - “Postcards, not ticket stubs” (Tasse et al., 2017)
- Find clever study designs or data comparisons, establish *panels*, etc.

Dan Tasse, Zichen Liu, Alex Sciuto, and Jason I. Hong. 2017. State of the geotags: Motivations and recent changes. *Proceedings of the Eleventh International AAAI Conference on Web and Social Media (ICWSM 2017)*, pp. 250-259.



Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

Platform effects in social media

Momin M. Malik and Jürgen Pfeffer. 2016. Identifying platform effects in social media data. In *Proceedings of the Tenth International AAAI Conference on Web and Social Media (ICWSM-16)*, pages 241-249. May 18-20, 2016, Cologne, Germany. Expanded version (2018):

https://www.mominmalik.com/malik_chapter2.pdf



Design can cause/change behavior

Introduction

Brief historical tour

Bias in geotagged tweets

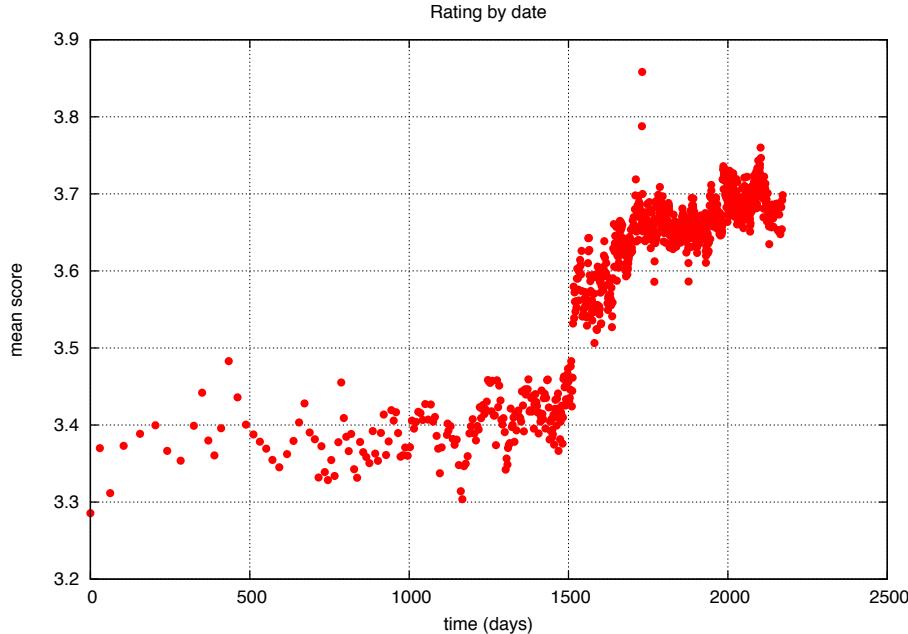
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Average Netflix movie ratings over time. Each point averages 100,000 rating instances.

Yehuda Koren. (2009). Collaborative filtering with temporal dynamics.



Social media platforms are businesses

Introduction

Brief historical tour

Bias in geotagged tweets

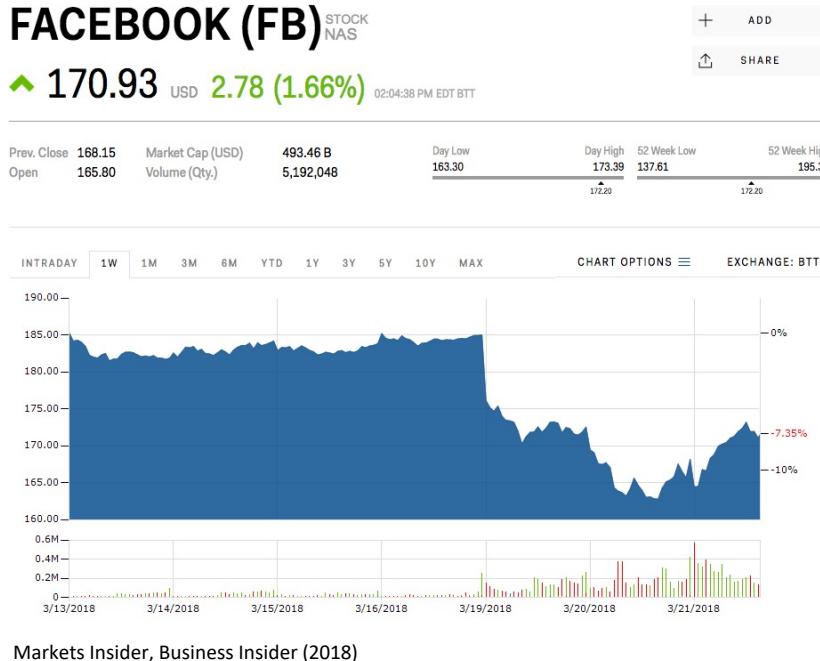
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Markets Insider, Business Insider (2018)

- Not neutral utilities or research environments
- Platform engineers try to shape user behavior towards desirable ends



Sites try to grow their users' networks

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

LinkedIn homepage featuring a banner: "It's easier than ever to grow your professional network". Below it, a section titled "INTRODUCING THE NEW" has a box around "People You May Know".

Twitter search results for "Who to follow". The first result is Keton Kakkar (@KetonKakkar), described as Afghan American / Child of Immigrants | @PhillipsAcademy / @Swarthmore | formerly @BKCHarvard | Editor @swatgazette. Followed by Frank Pasquale and monicabulger.

The second result is William Bumpas (@wwbumpas), described as Now in DC, prev @oioxford. Likes data, ethnography, tech, policy, media, critical theory, China, rural US, subversive memes. Loves any combo thereof. he/she/they. Followed by Prof Gina Neff and Oxford Internet Institute.

The third result is Rich Boroff (@boroff), described as Running (a minor part of) the computing infrastructure for a major university in the Boston, MA area, and trying to keep the bad guys at bay. Followed by Berkman Klein Center for Internet & Society.



Recommending “friend-of-a-friend”

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

Search Facebook Dann Home

People you may know

	Sara Anderson Severance Denver, Colorado Rachelle Albright and 10 other mutual friends	Add Friend Remove
	Anne Walker (Anne Anderson) Sarah Frederick and 6 other mutual friends	Add Friend Remove
	Paul Dube Ryan Dube is a mutual friend.	Add Friend Remove
	Mark Rieder Lord Beaverbrook High School Justin Pot is a mutual friend.	Add Friend Remove

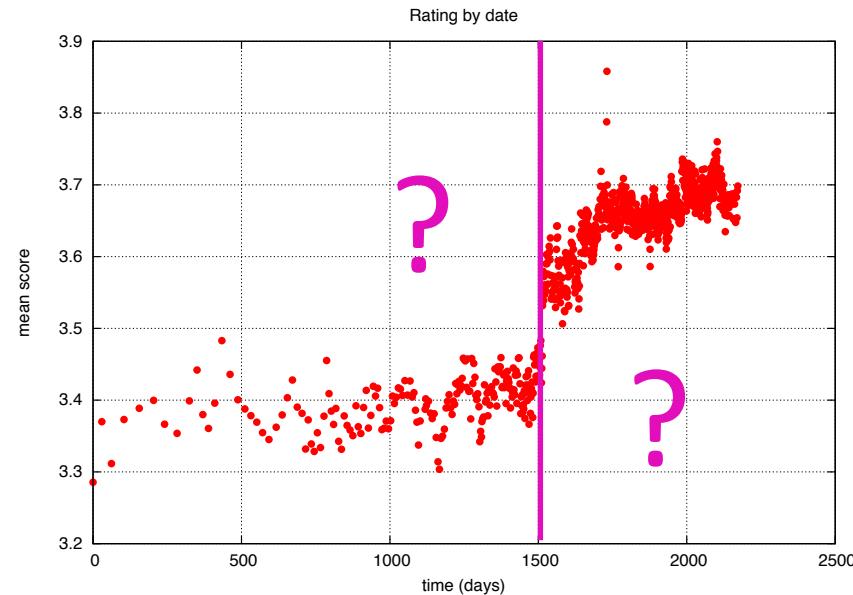
Dann Albright, makeuseof.com

Search for Friends
Find friends from all
Name
Search for someone
Home Town
 Prescott, Wisconsin
Enter another city
Current location
 Denver, Colorado
Enter another city
High School
 Prescott High School
Enter another high school



Behavior, or platform effects?

- When we measure behavior, what are we really measuring?
People's behavior, or platform effects?
- How, as outsiders, can we find out?



Average Netflix movie ratings over time. Each point averages 100,000 rating instances.

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Data artifacts can reveal inner workings

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References





Data artifacts as natural experiments

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

- Regression Discontinuity (RD) Design (technically, Interrupted Time Series, ITS) estimates causality

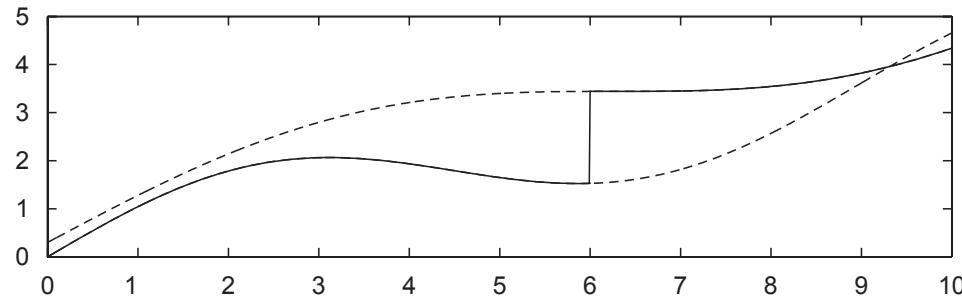


Fig. 2 from Imbens and Lemieux (2008): Potential and observed outcome regression functions.

- The difference between “before” and “after” estimates the *local average treatment effect*

Guido W. Imbens and Thomas Lemieux. (2008). Regression discontinuity designs: A guide to practice. *Journal of Econometrics*, 142(2): 615–635.



Case: Facebook's "People You May Know"

Introduction

Brief historical tour

Bias in geotagged tweets

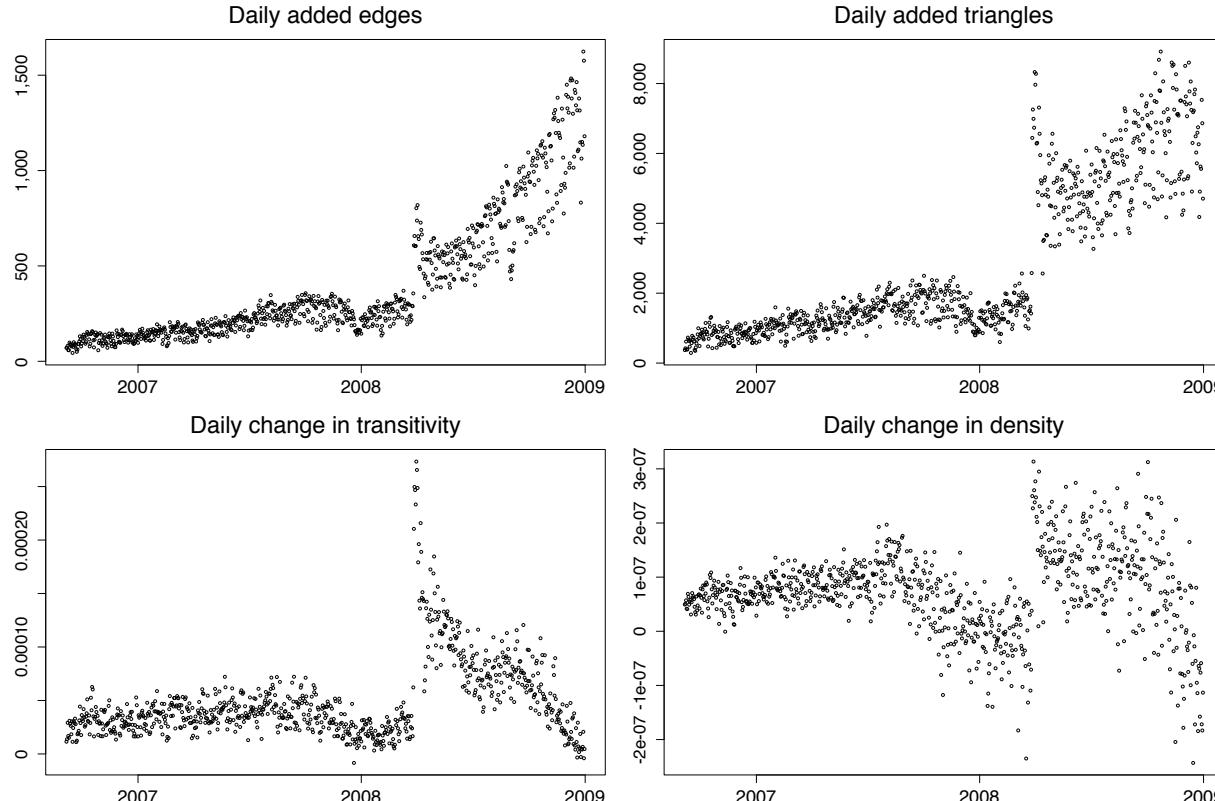
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References





PYMK changed the Facebook network!

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

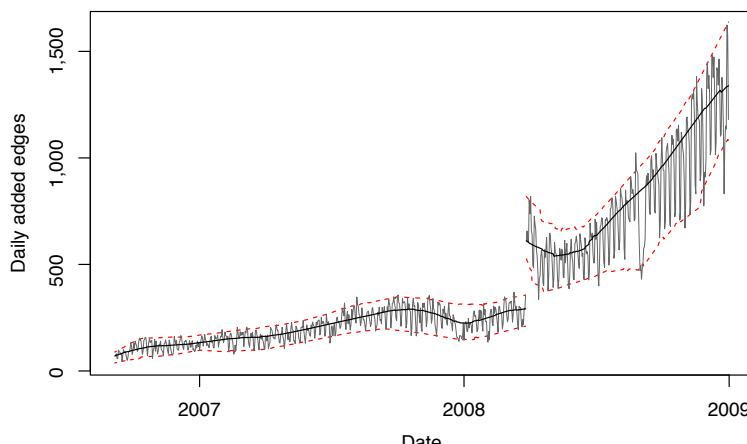
Hierarchy of limitations in machine learning

Problems of cross-validation

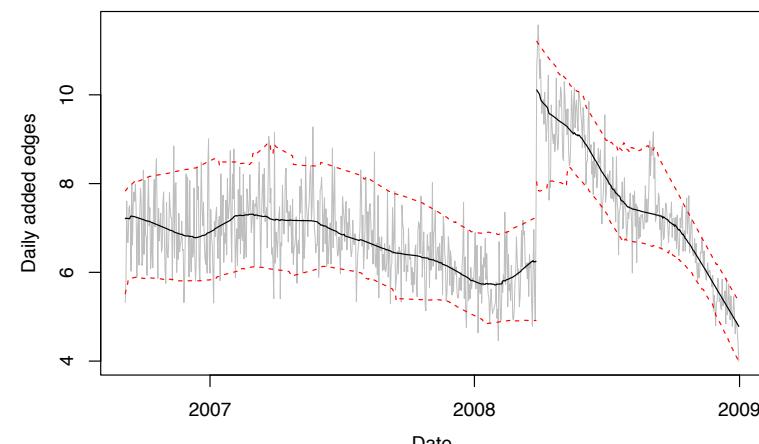
Summary and conclusion

References

- Facebook links: +300 new edges per day (x2)



- Triangles: +3.8 triangles per edge (x1.62)





Responses to platform effects

- Investigate: how do Facebook “friendship” fail to generalize? What about the Facebook social network?
- Platform effects are phenomena to study in themselves!
- Data artifacts as natural experiments

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

Hierarchy of limitations in machine learning

Momin M. Malik. 2020. A hierarchy of limitations in machine learning.

<https://arxiv.org/abs/2002.05193>

Data well-considered; models, not so much

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

Critical perspective on measurement in digital trace data and machine learning

34 of 75

frontiers in Big Data

REVIEW
published: 11 July 2019
doi: 10.3389/fdata.2019.00019

Accepted: 27 May 2019
Published: 11 July 2019
Citation:

Social Data: Biases, Methodological Pitfalls, and Ethical Boundaries

Alexandra Oiteanu^{1,*}, Carlos Castillo², Fernando Diaz² and Emre Kiciman¹

¹ Microsoft Research, New York, NY, United States; ² Microsoft Research, Montreal, QC, Canada; ³ Department of Information and Communication Technologies, Universitat Pompeu Fabra, Barcelona, Spain; ⁴ Microsoft Research, Redmond, WA, United States

Social data in digital form—including user-generated content, expressed or implicit relations between people, and behavioral traces—are at the core of popular applications and platforms, driving the research agenda of many researchers. The promises of social data are many, including understanding “what the world thinks” about a social issue, brand, celebrity, or other entity, as well as enabling better decision-making in a variety of fields including public policy, healthcare, and economics. Many academics and practitioners have warned against the naïve usage of social data. There are biases and inaccuracies occurring at the source of the data, but also introduced during processing. There are methodological limitations and pitfalls, as well as ethical boundaries and unexpected consequences that are often overlooked. This paper recognizes the rigor with which these issues are addressed by different researchers varies across a wide range. We identify a variety of menaces in the practices around social data use, and organize them in a framework that helps to identify them.

OPEN ACCESS

Edited by:
Juergen Pfeffer,
University of Bayreuth, Germany

Reviewed by:
Kenneth Joseph,
University of Buffalo, United States;
Mark M. Mouloua,
Harvard University, United States

***Correspondence:**
alexandra.oiteanu@microsoft.com

Specialty section:
This article was submitted to
Data Mining and Management,
a section of the journal
Frontiers in Data Science

Received: 26 February 2019
Accepted: 27 May 2019
Published: 11 July 2019
Citation:

For your own sanity, you have to remember that not all problems can be solved. Not all problems can be solved, but all problems can be illuminated. – Ursula Franklin¹

Keywords: social media, user data, biases, evaluation, ethics

1. INTRODUCTION

We use *social data* as an umbrella concept for all kind of digital traces produced by or about users, with an emphasis on content explicitly written with the intent of communicating or interacting with others. Social data typically comes from *social software*, which provides an intermediary or a focus for a social relationship or interaction. It includes a variety of platforms—like for social media and social networking (e.g., Facebook), question and answering (e.g., Quora), or collaboration (e.g., Wikipedia)—and purposes from finding information (Whitelock, 2013) to keeping in touch with friends (Lampe et al., 2008). Social software enables the social web, a class of websites “in which user participation is the primary driver of value” (Gruber, 2008).

The social web enables access to social traces at a scale and level of detail, both in breadth and depth, impractical with conventional data collection techniques, like surveys or user



Approaches to research

Introduction

Brief historical tour

Bias in geotagged tweets

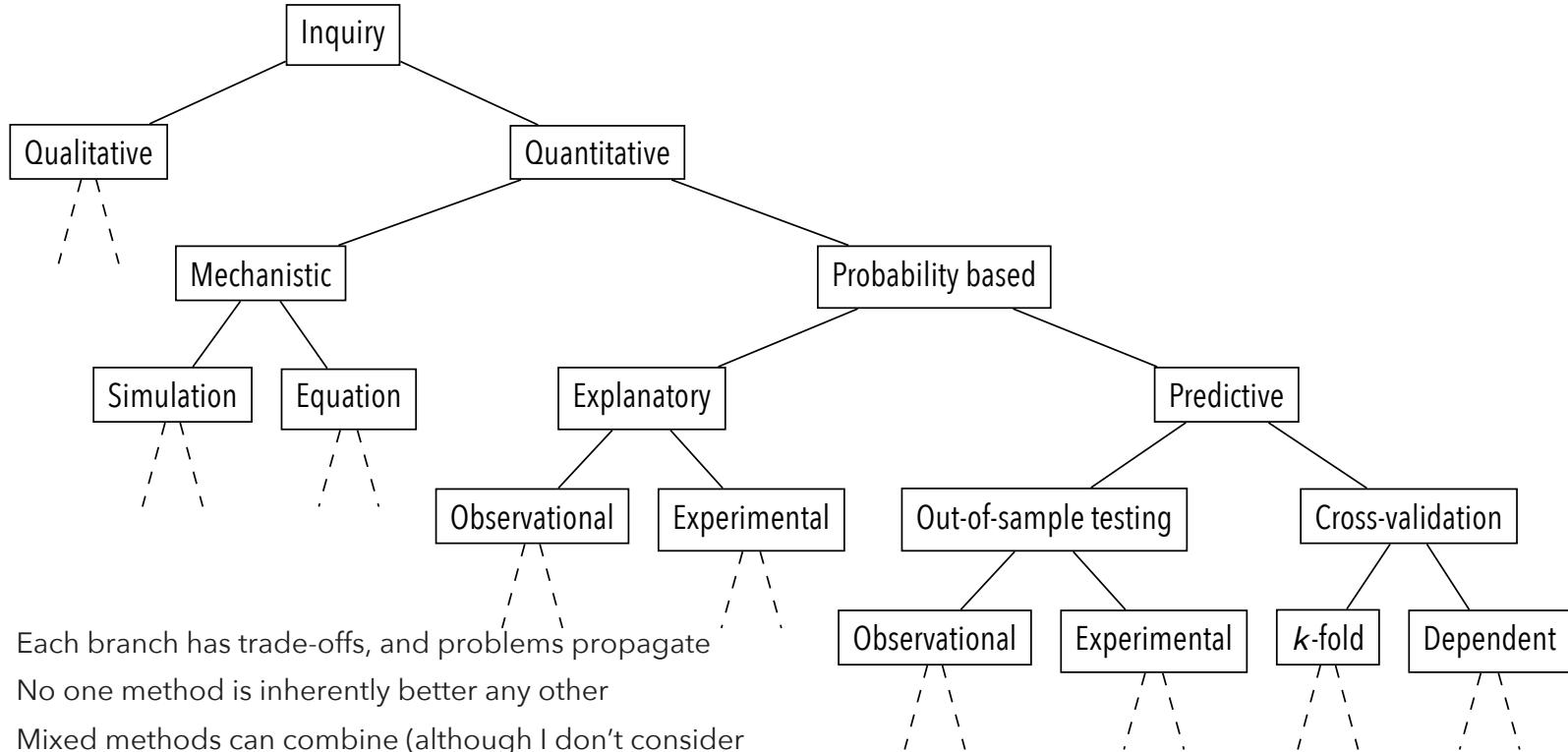
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References





Quantification locks in meaning

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

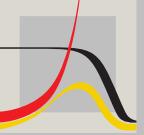
Problems of cross-validation

Summary and conclusion

References



- Qualitative research can get directly at how things are multifaceted, heterogeneous, intersubjective
- Quantification/ measurements lock in one meaning; and frequently are *proxies*, which are imperfect



Challenges of quantification/ measurement

Introduction

Brief historical tour

Bias in
geotagged
tweets

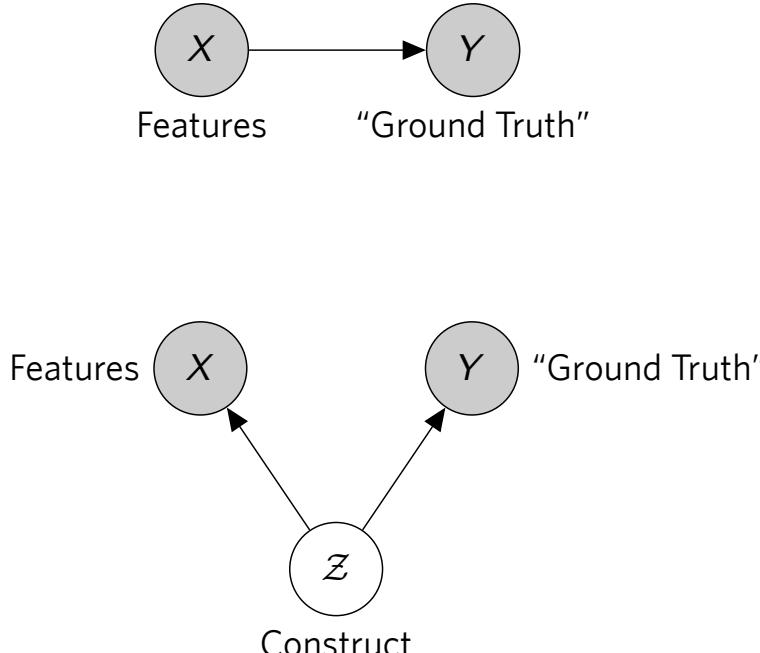
Platform
effects

**Hierarchy of
limitations in
machine
learning**

Problems of
cross-
validation

Summary and
conclusion

References



- *Constructs*: primitives of social science
 - What we care about
 - Often unobservable (and hypothetical/subjective, e.g. friendship)
 - Proxies always give errors (for binary constructs: false negatives and false positives)
 - E.g., Google maps usage is not traffic



Constructs: Subjective, multifaceted

Introduction

Brief historical tour

Bias in geotagged tweets

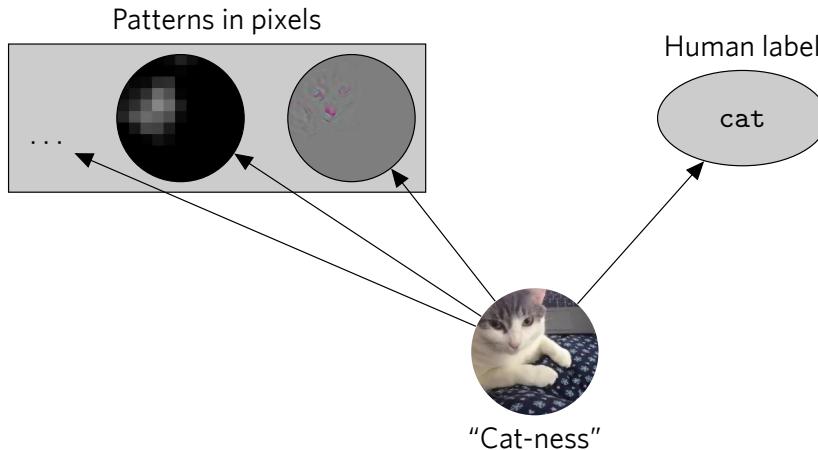
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

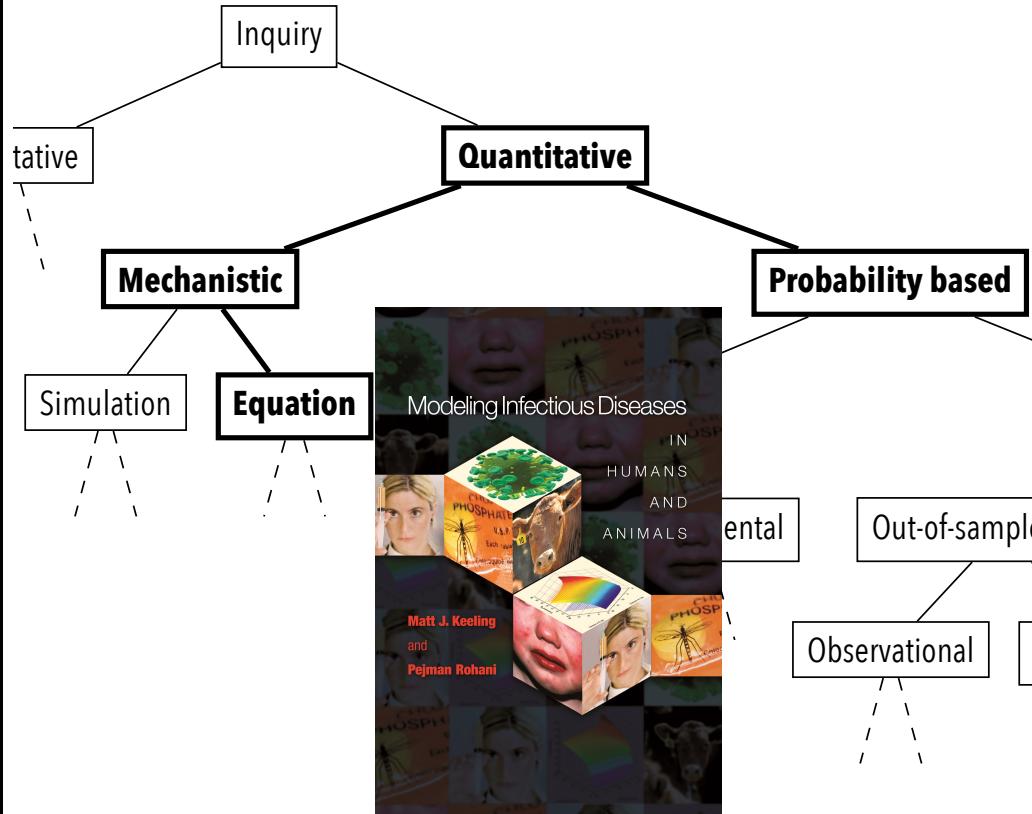
Summary and conclusion

References





Stats and ML use central tendencies



- Statistics and machine learning are the only options to both directly use data and account for variability
 - They do so via central tendency
 - This requires multiple observations, and independence assumptions



Stats and ML use central tendencies

Introduction

Brief historical tour

Bias in geotagged tweets

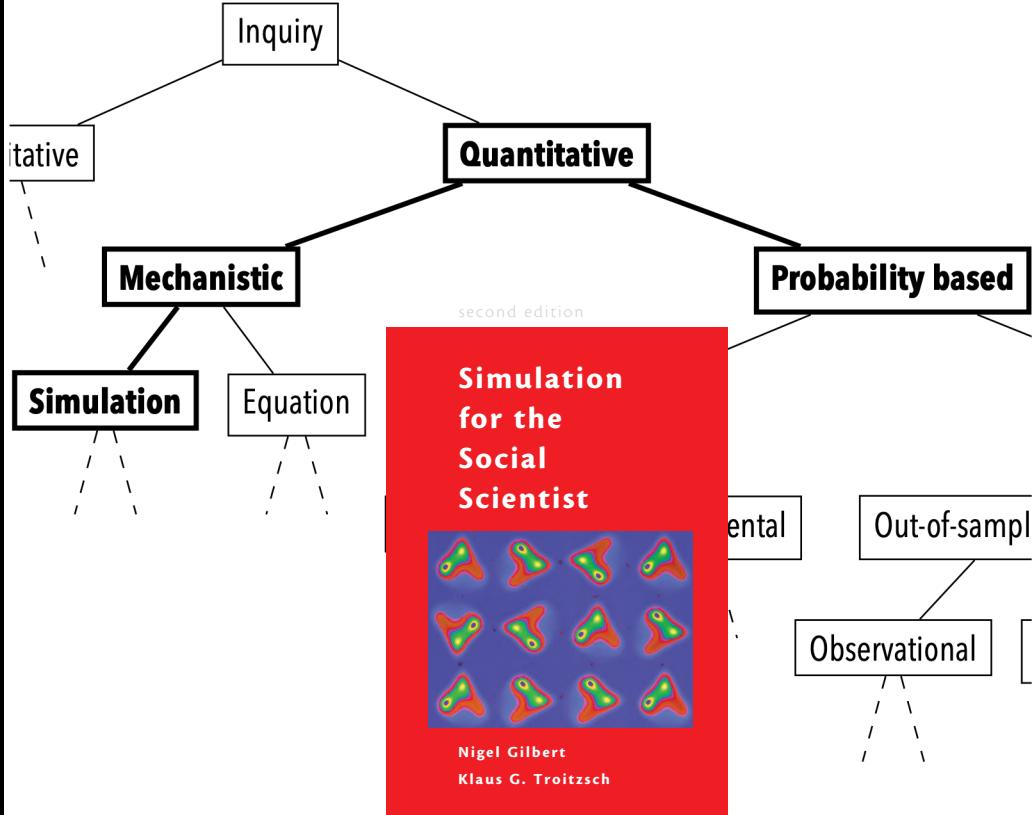
Platform effects

Hierarchy of limitations in machine learning

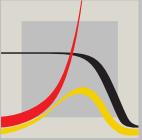
Problems of cross-validation

Summary and conclusion

References

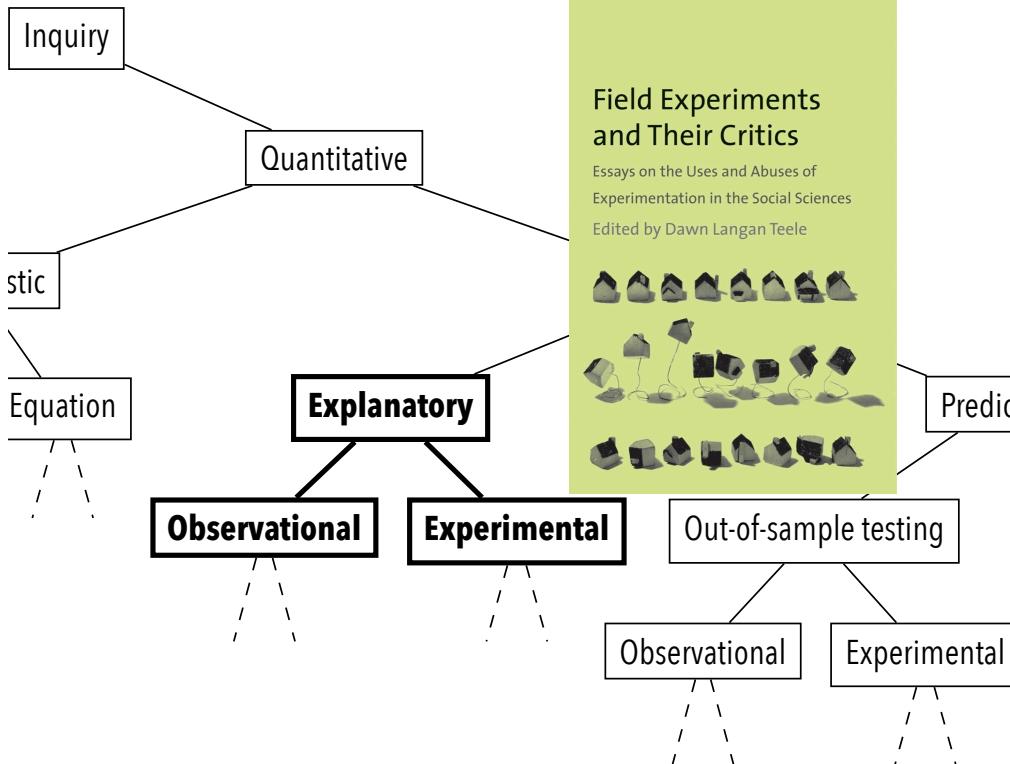


- (Statistics uses *numerical* simulations, and simulation modeling uses statistical *summaries*, but they are distinct types of models)
- (Agent-based simulation also ends up using central tendencies to summarize a response surface)
- (ABMs generally cannot be used for prediction, are only appropriate when we can't do statistics)



Causality is hard, maybe too hard

Introduction
Brief historical tour
Bias in geotagged tweets
Platform effects
Hierarchy of limitations in machine learning
Problems of cross-validation
Summary and conclusion
References

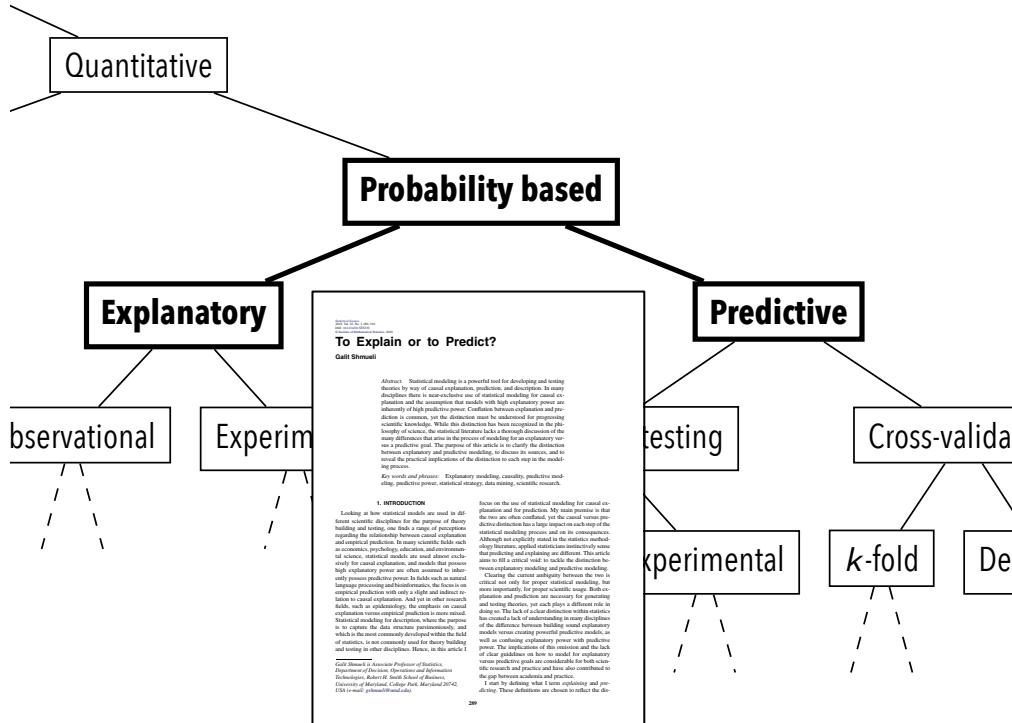


- Properly controlled experiments lack ecological validity
- Observational inference can never totally account for the possibility of hidden confounders, which can frustrate even the most perfect application of causal techniques (Arceneaux, Gerber, & Green, 2013)



ML is “prediction” only

- “Predictions” are defined as what minimizes loss
 - I.e., *correlations*
 - Non-causal correlations can sometimes predict well, but they frequently don’t explain, and can fail unexpectedly





Defining machine learning

Introduction

Brief historical tour

Bias in geotagged tweets

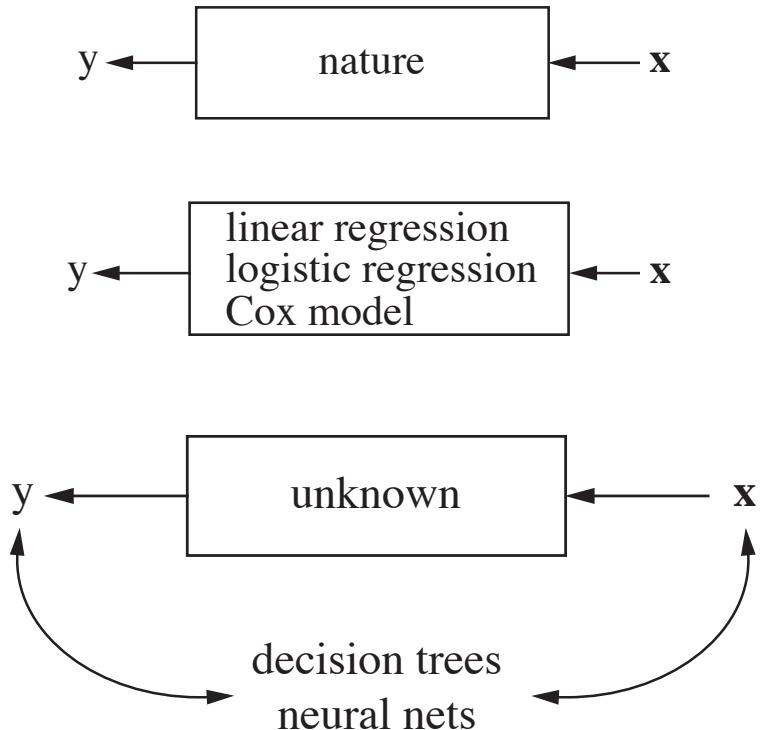
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



- Machine learning: An instrumental use of correlations to try and *mimic* the outputs of a target system (rather than trying to understand causal relationships between inputs and outputs). Focus on highly flexible “curve-fitting” methods. (Diagram: Breiman, 2001. See also Jones, 2018)



Why are these different goals?

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

$$\hat{y}$$

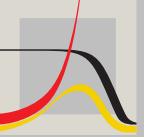
$$\hat{\beta}$$

Spurious (non-causal) correlations may fit robustly

- Breiman 2001: Prediction problems
- Shmueli 2010: To predict
- Kleinberg et al. 2015: “Umbrella problems”
- Mullainathan and Spiess 2017: \hat{y}

Carefully built models that capture causality (or “pure” associations) may fit poorly overall

- Breiman 2001: Information
- Shmueli 2010: To explain
- Kleinberg et al. 2015: “Rain dance problems”
- Mullainathan and Spiess 2017: $\hat{\beta}$



ML: Only external validity

Introduction

Brief historical tour

Bias in geotagged tweets

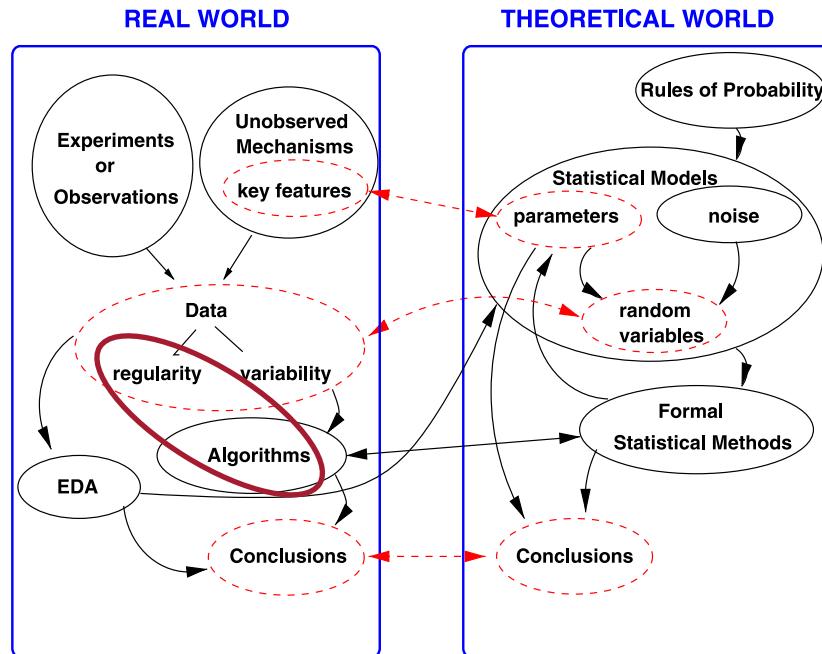
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Kass, 2011, Stat. Sci.

Adapted from Borgatti, 2012



Not an obvious usage of “predict”

Introduction

Brief historical tour

Bias in geotagged tweets

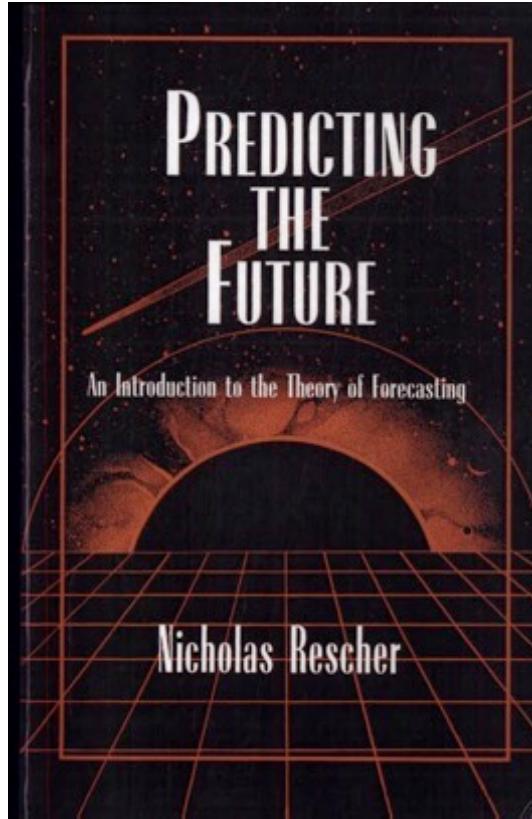
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



88 ■ PREDICTING THE FUTURE

TABLE 6.1: A SURVEY OF PREDICTIVE APPROACHES

Predictive Approaches	Linking Mechanism	Methodology Of Linkage
UNFORMALIZED/JUDGMENTAL		
judgmental estimation	expert informants	informed judgment
FORMALIZED/INFERENTIAL		
RUDIMENTARY (ELEMENTARY)		
trend projection	prevailing trends	projection of prevailing trends
curve fitting	geometric patterns	subsumption under an established pattern
circumstantial analogy	comparability groupings	assimilation to an analogous situation
SCIENTIFIC (SOPHISTICATED)		
indicator coordination	causal correlations	statistical subsumption into a correlation
law derivation (nomic)	accepted laws (deterministic or statistical)	inference from accepted laws
phenomenological modeling (analogical)	formal models (physical or mathematical)	analogizing of actual (“real-world”) processes with presumably isomorphic model process



Can't *intervene* based on correlations

Introduction

Brief historical tour

Bias in geotagged tweets

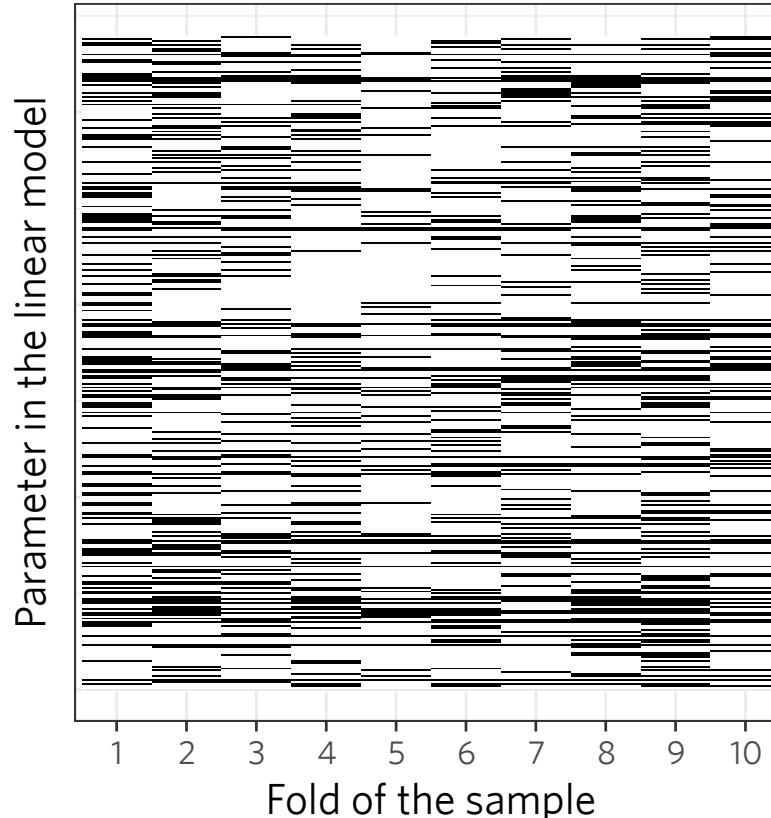
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

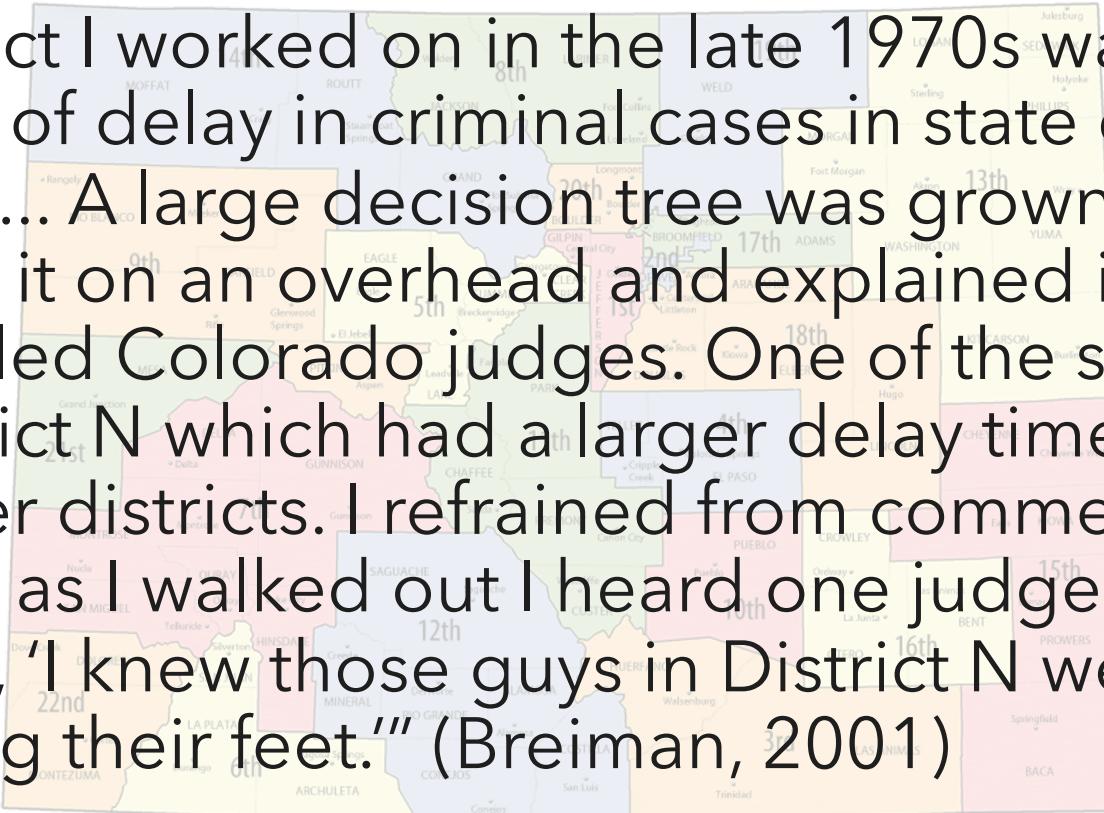


- Very different sets of correlations can “predict” (fit) equally well (Mullainathan and Spiess 2017)
 - Breiman (2001) called this the “Rashomon Effect”
- But different fits suggest very different interventions



Interpretability: A red herring?

"A project I worked on in the late 1970s was the analysis of delay in criminal cases in state court systems... A large decision tree was grown, and I showed it on an overhead and explained it to the assembled Colorado judges. One of the splits was on District N which had a larger delay time than the other districts. I refrained from commenting on this. But as I walked out I heard one judge say to another, 'I knew those guys in District N were dragging their feet.'" (Breiman, 2001)





Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

Problems of cross validation



ML performance claims are from cross-validation

Introduction

Brief historical tour

Bias in geotagged tweets

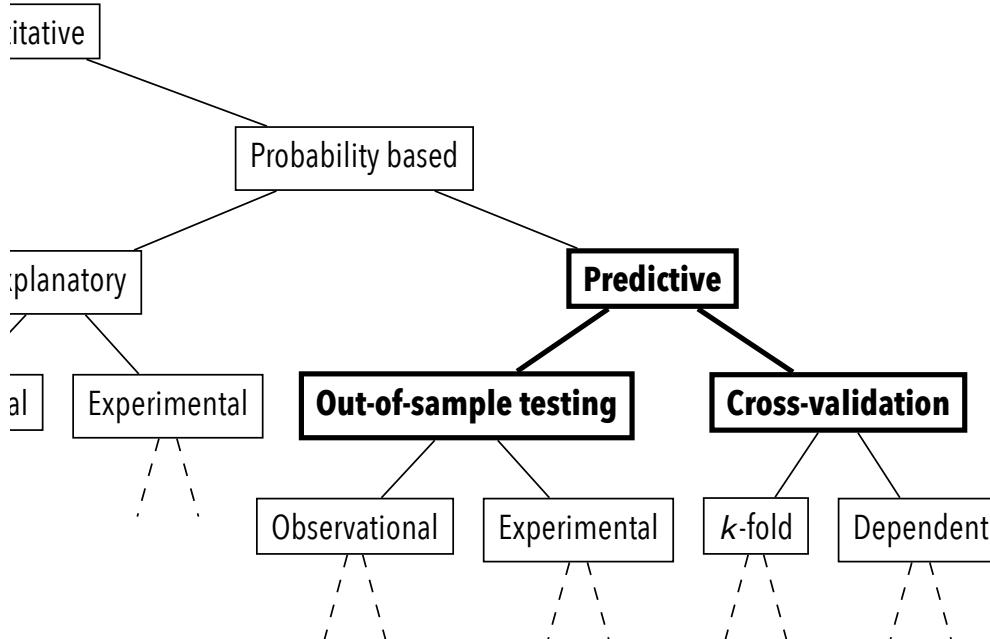
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



- Rescher (1998) notes every prediction involves a meta-prediction: predict whether the prediction works
- Cross-validation is meta-prediction for ML
- But, how well does cross-validation work?
 - "Professor Breiman emphasizes the importance of performance on the test sample. However, this can be overdone. The test sample is supposed to represent the population to be encountered in the future. But in reality, it is usually a random sample of the current population. High performance on the test sample does not guarantee high performance on future samples, **things do change.**" (Hoadley 2001, discussant of Breiman)



Purpose of cross-validation

Introduction

Brief historical tour

Bias in geotagged tweets

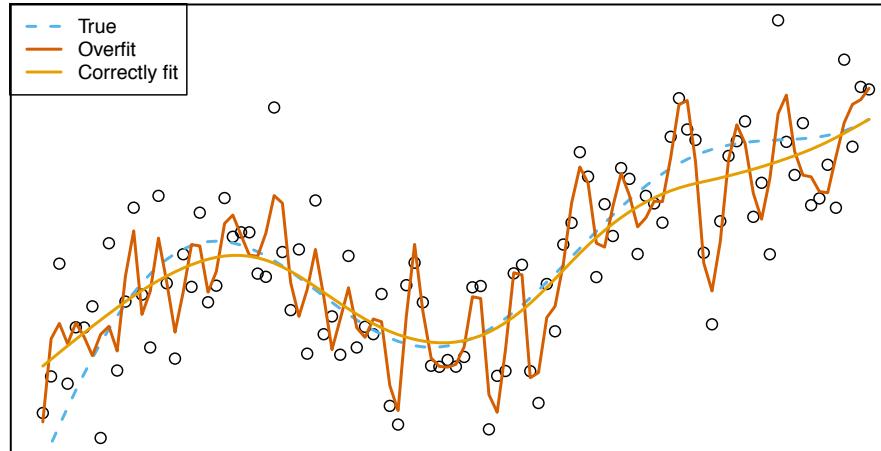
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



- If we are no longer guided by theory, and use automatic methods, we risk overfitting: fitting to the noise, not the data



Intuition for cross-validation

Introduction

Brief historical tour

Bias in geotagged tweets

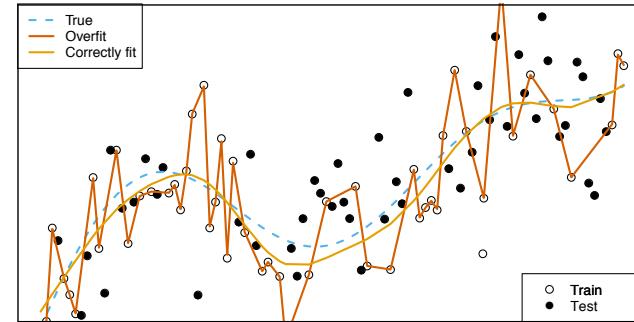
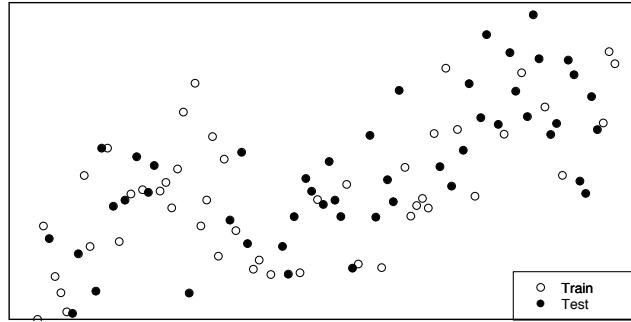
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



- Idea: if we split data into two parts, the signal should be the same but the noise would be different
- *Cross validation*: Fitting the model on one part of the data, and “testing” on the other



Overfitting on the test set

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

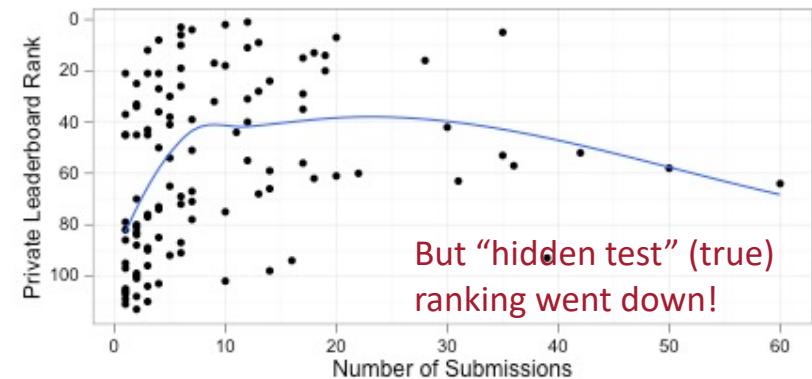
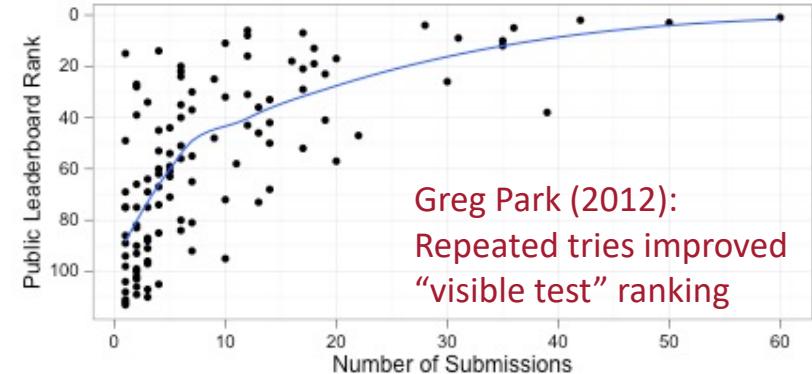
Hierarchy of limitations in machine learning

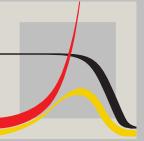
Problems of cross-validation

Summary and conclusion

References

- Re-using a test set can overfit! (Dwork et al., 2015)
 - “in industry and academia, there is sometimes a little tinkering, which involves peeking at the test sample. The result is some bias in the test sample or cross-validation results. This is the same kind of tinkering that upsets test of fit pureness.” (Hoadley 2001, discussant of Breiman)
- Happens in Kaggle, which has public leaderboard (visible throughout) and private leaderboard (revealed only at end of competition)





Problems of dependencies

Introduction

Brief historical tour

Bias in geotagged tweets

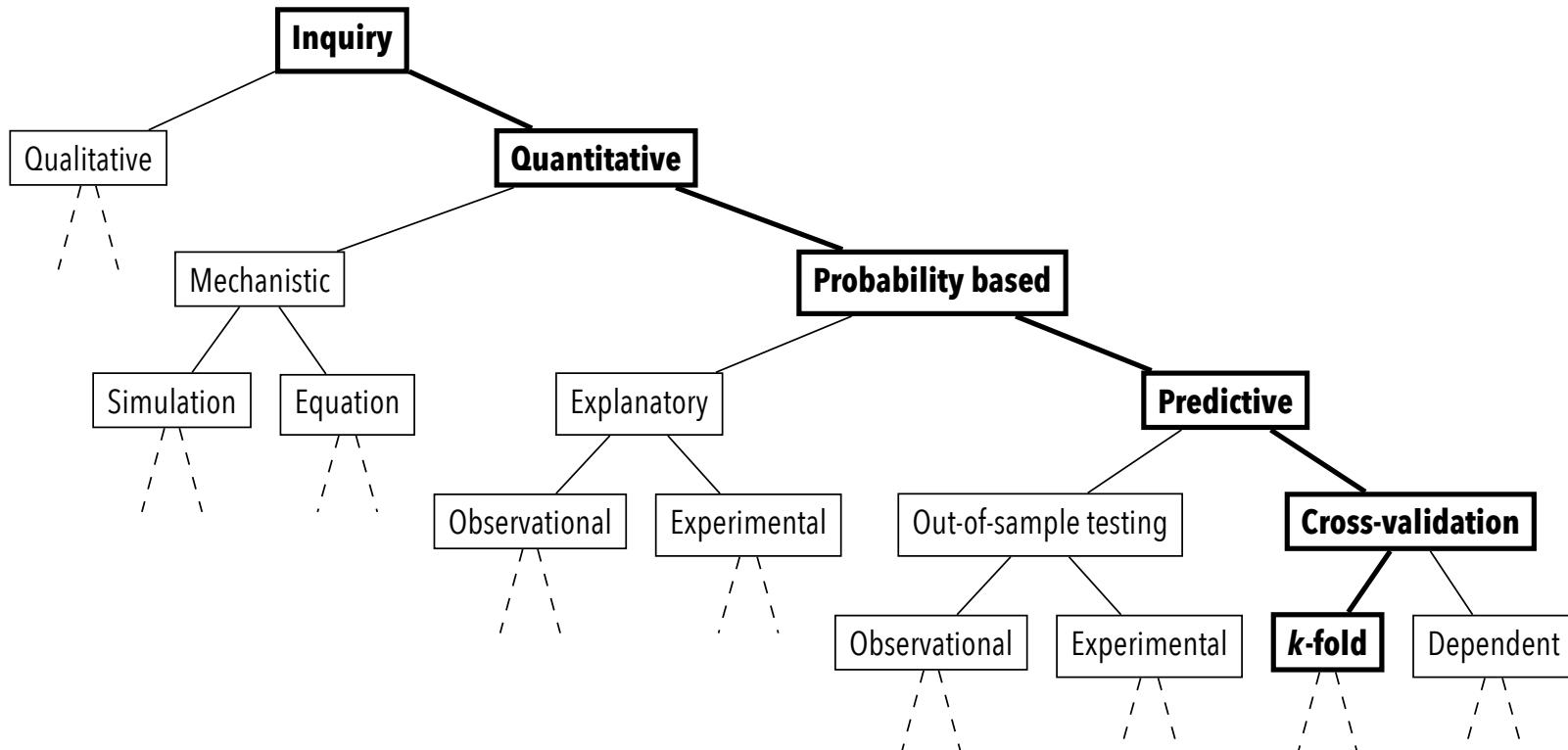
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References





Classic argument for CV

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

$$\begin{aligned}\text{Err}(\hat{\mu}) &= \frac{1}{n} \mathbb{E}_f \|Y^* - \hat{Y}\|_2^2 \\&= \frac{1}{n} \left[\mathbb{E}_f \|Y^*\|_2^2 + \mathbb{E}_f \|\hat{Y}\|_2^2 - 2\mathbb{E}_f (Y^{*T} \hat{Y}) \right] \\&= \frac{1}{n} \left[\mathbb{E}_f \|Y^*\|_2^2 + \mathbb{E}_f \|\hat{Y}\|_2^2 - 2 \text{tr } \mathbb{E}_f (Y^* \hat{Y}^T) \right] \\&\quad + \frac{1}{n} \left[-\mu^T \mu + \mathbb{E}_f (\hat{Y})^T \mathbb{E}_f (\hat{Y}) + 2 \text{tr } \mu \mathbb{E}_f (\hat{Y})^T \right] \\&\quad + \frac{1}{n} \left[-\mu^T \mu - \mathbb{E}_f (\hat{Y}) \mathbb{E}_f (\hat{Y})^T - 2\mu^T \mathbb{E}_f (\hat{Y}) \right] \\&= \frac{1}{n} \left[\text{tr } \Sigma + \|\mu - \mathbb{E}(\hat{Y})\|_2^2 + \text{tr } \text{Var}_f(\hat{Y}) - 2 \text{tr } \text{Cov}_f(Y^*, \hat{Y}) \right] \\&= \text{irreducible error} + \text{bias}^2 + \text{variance} - \text{optimism}\end{aligned}$$



Apply this to non-iid data

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

- Imagine we have, for $\Sigma_{ii} = \sigma^2$ and $\Sigma_{ij} = \rho\sigma^2$, $i \neq j$

$$\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{X} \end{bmatrix} \beta, \begin{bmatrix} \Sigma & \rho\sigma^2 \mathbf{1} \mathbf{1}^T \\ \rho\sigma^2 \mathbf{1} \mathbf{1}^T & \Sigma \end{bmatrix} \right)$$

- Then, optimism in the training set is:

$$\frac{2}{n} \operatorname{tr} \operatorname{Cov}_f(Y_1, \hat{Y}_1) = \frac{2}{n} \operatorname{tr} \operatorname{Cov}_f(Y_1, \mathbf{H}Y_1) = \frac{2}{n} \operatorname{tr} \mathbf{H} \operatorname{Var}_f(Y_1) = \frac{2}{n} \operatorname{tr} \mathbf{H} \Sigma$$

- But test set also has nonzero optimism!

$$\frac{2}{n} \operatorname{tr} \operatorname{Cov}_f(Y_2, \hat{Y}_1) = \frac{2}{n} \operatorname{tr} \operatorname{Cov}_f(Y_2, \mathbf{H}Y_1) = \frac{2\rho\sigma^2}{n} \operatorname{tr} \mathbf{H} \mathbf{1} \mathbf{1}^T = 2\rho\sigma^2$$



Simulating the toy example

Introduction

Brief historical tour

Bias in geotagged tweets

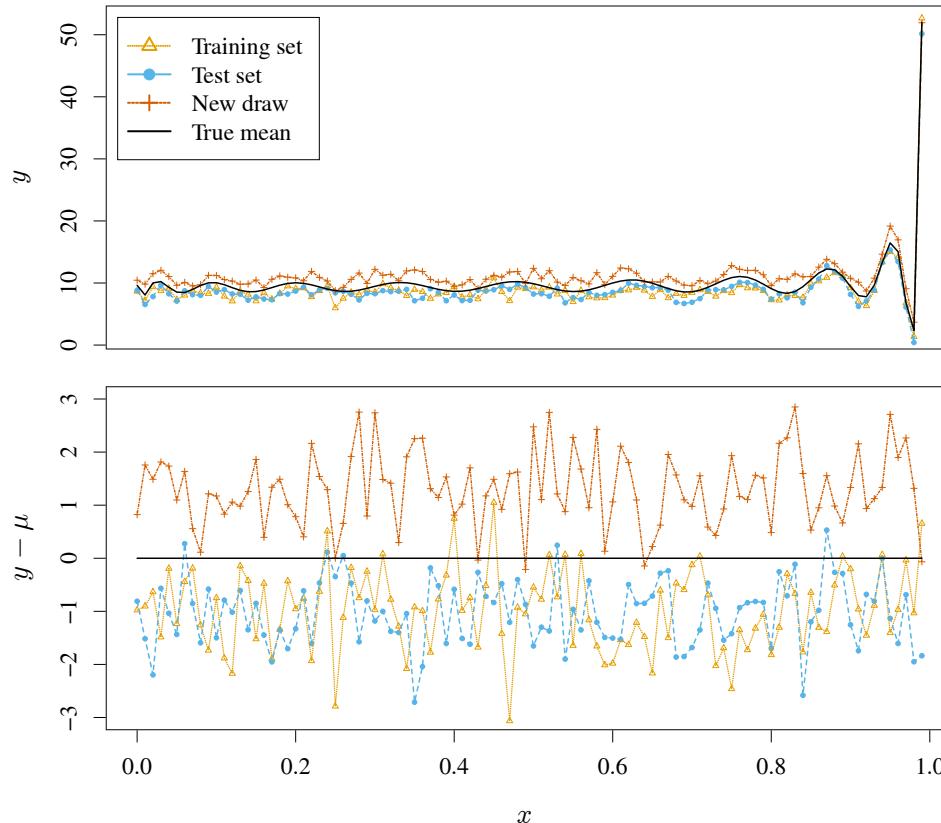
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

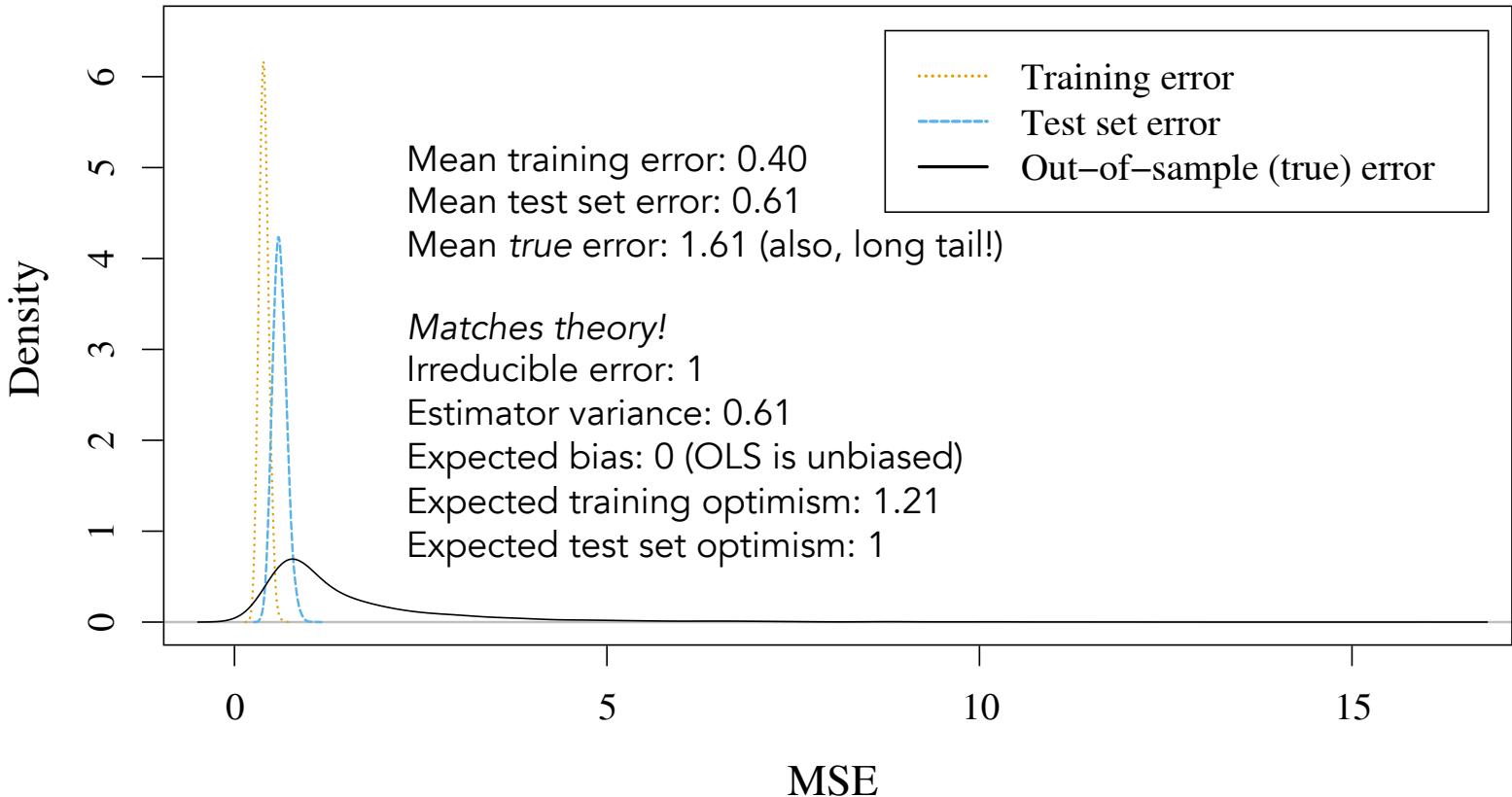
References





Out-of-sample MSE: *much worse!*

- Introduction
- Brief historical tour
- Bias in geotagged tweets
- Platform effects
- Hierarchy of limitations in machine learning
- Problems of cross-validation**
- Summary and conclusion
- References





Many real-world examples

- There are indeed cases where cross-validation assessments of machine learning performance fail!
- Time series: do cross-validation in blocks
 - Otherwise, “time traveling,” gives great performance
- Activity recognition: “leave one subject out” cross validation performs far worse (i.e., more honestly)
- Necessary but not sufficient; underlying causal processes can introduce unobserved variance, destroying previously-holding correlations

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Application to networks

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

	Y	X_1	X_2	\dots	X_d
1	y_1	x_{11}	x_{12}	\dots	x_{1d}
2	y_2	x_{21}	x_{22}	\dots	x_{2d}
\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
n	y_n	x_{n1}	x_{n2}	\dots	x_{nd}



$index$	$from$	to	Y	W_1	W_2	W_3	\dots
e_1	1	2	y_{12}	$\mathbf{1}(x_{11} = x_{21})$	$x_{12} - x_{22}$	x_{13}	\dots
e_2	2	3	y_{23}	$\mathbf{1}(x_{11} = x_{31})$	$x_{12} - x_{32}$	x_{13}	\dots
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
e_{n+1}	2	1	y_{21}	$\mathbf{1}(x_{21} = x_{11})$	$x_{22} - x_{12}$	x_{23}	\dots
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$e_{2\binom{n}{2}}$	$n-1$	n	$y_{(n-1)n}$	$\mathbf{1}(x_{(n-1)1} = x_{n1})$	$x_{(n-1)2} - x_{n2}$	$x_{(n-1)3}$	\dots



But dyads are dependent too!

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

Factor graph	Parameter name	Network Motif	Parameterization	Matrix notation
	-mutual dyads		$\sum_{i < j} A_{ij} A_{ji}$	$\frac{1}{2} \text{tr}(\mathbf{AA}^T)$
	-in-two-stars		$\sum_{(i,j,k)} A_{ji} A_{ki}$	$\text{sum}(\mathbf{AA}^T) - \text{tr}(\mathbf{AA}^T)$
	-out-two-stars		$\sum_{(i,j,k)} A_{ij} A_{ik}$	$\text{sum}(\mathbf{A}^T \mathbf{A}) - \text{tr}(\mathbf{A}^T \mathbf{A})$
	geom. weighted out-degrees	—	$\sum_i \exp\{-\alpha \sum_k A_{ik}\}$	$\text{sum}(\exp\{-\alpha \text{rowsum}(\mathbf{A})\})$
	geom. weighted in-degrees	—	$\sum_j \exp\{-\alpha \sum_k A_{kj}\}$	$\text{sum}(\exp\{-\alpha \text{colsum}(\mathbf{A})\})$
	-alternating transitive k-triplets		$\lambda \sum_{i,j} A_{ij} \left\{ 1 - \left(1 - \frac{1}{\lambda}\right) \sum_{k \neq i,j} A_{ik} A_{kj} \right\}$	$\lambda \sum (\mathbf{A}^{(\cdot)} \left(1 - \left(1 - \frac{1}{\lambda}\right) \mathbf{AA} - \text{diag}(\mathbf{AA}) \right))$
	-alternating indep. two-paths		$\lambda \sum_{i,j} \left\{ 1 - \left(1 - \frac{1}{\lambda}\right) \sum_{k \neq i,j} A_{ik} A_{kj} \right\}$	$\lambda \sum \left(1 - \left(1 - \frac{1}{\lambda}\right) \mathbf{AA} - \text{diag}(\mathbf{AA}) \right)$
	-two-paths (mixed two-stars)	—	$\sum_{(i,k,j)} A_{ij} A_{kj}$	$\text{sum}(\mathbf{AA}) - \text{tr}(\mathbf{AA})$
	-transitive triads		$\sum_{(i,j,k)} A_{ij} A_{jk} A_{ik}$	$\text{tr}(\mathbf{AAA}^T)$
	-activity effect		$\sum_i X_i \sum_j A_{ij}$	$\text{sum}(\mathbf{X}^{(\cdot)} \text{rowsum}(\mathbf{A}))$
	-popularity effect		$\sum_j X_j \sum_i A_{ij}$	$\text{sum}(\mathbf{X}^{(\cdot)} \text{colsum}(\mathbf{A}))$
	-similarity effect		$\sum_{i,j} A_{ij} \left(1 - \frac{ X_i - X_j }{\max_{k,l} X_k - X_l } \right)$	$\text{sum}(\mathbf{A}^{(\cdot)} \mathbf{S})$

Graphical model and matrix notations for ERGM specification terms given in: Snijders et al. 2006. Joint work with Antonis Manousis and Naji Shajarisales, 2018.



Covariance structure of edges ($n = 15$)

Introduction

Brief historical tour

Bias in geotagged tweets

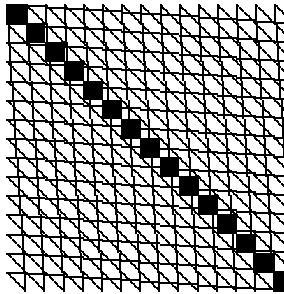
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

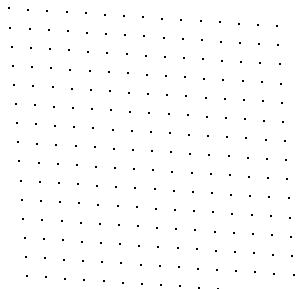
Summary and conclusion

References

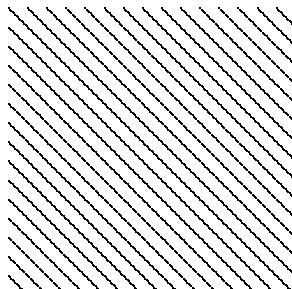


Total covariance between dyads

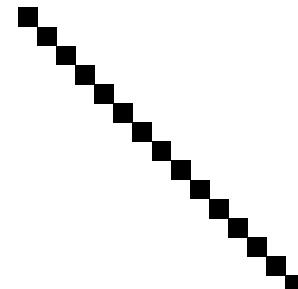
- The pairs of edges that are present together, or aren't present together
- Note: A theoretical construct, since we only see edges once (or once per time slice)



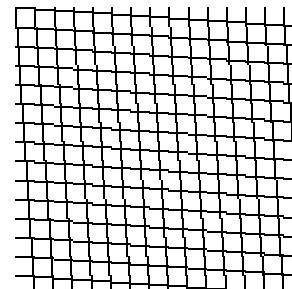
Mutual dyads



In-2-stars



Out-2-stars



2-paths



So, what to do?

Introduction

Brief historical tour

Bias in geotagged tweets

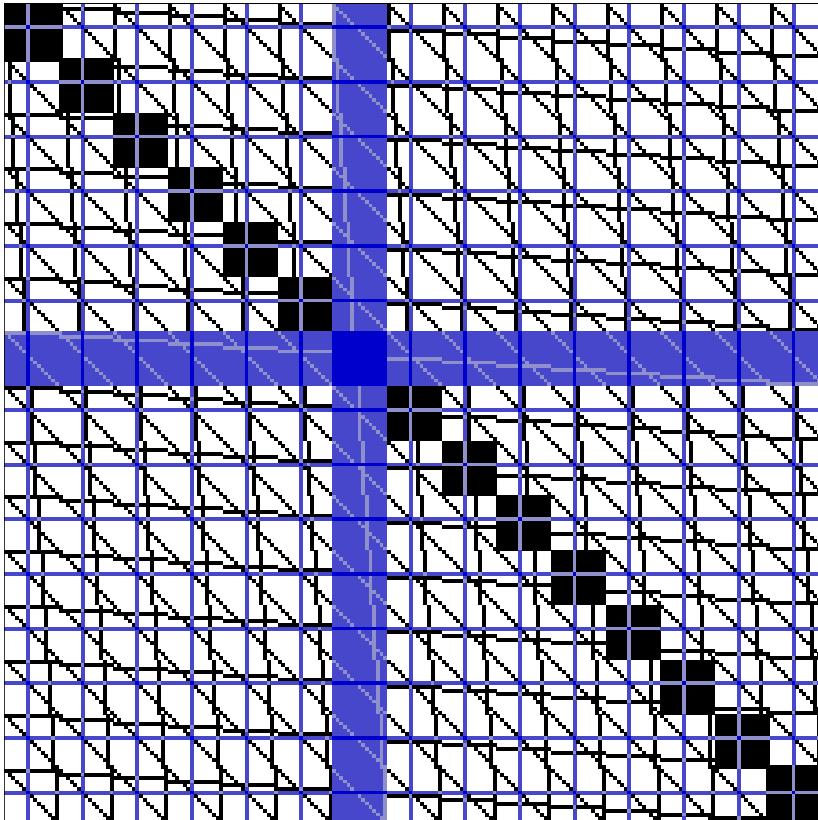
Platform effects

Hierarchy of limitations in machine learning

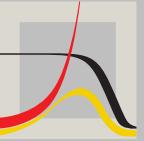
Problems of cross-validation

Summary and conclusion

References



- Partition nodes into training and test sets?
 - Breaks up triads; omitted edges “share” information across training and test (diagram: blue are edges that include node 7)
- Partition dyads?
 - Breaks up nodes; even worse
- Can’t *eliminate*, but can *minimize* optimism by careful data splitting



Importance of out-of-sample testing

Introduction

Brief historical tour

Bias in geotagged tweets

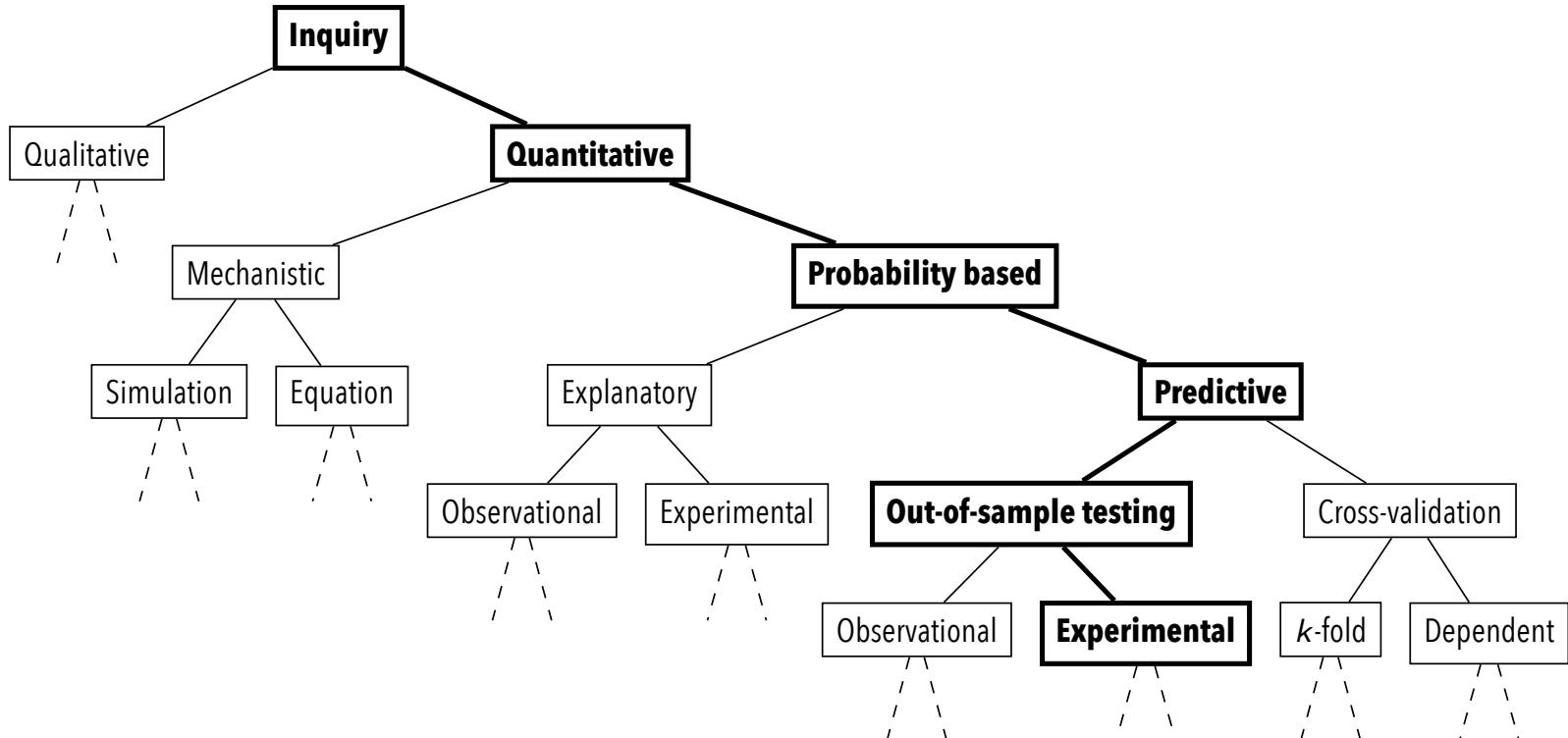
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References





“Things do change”

Introduction

Brief historical tour

Bias in geotagged tweets

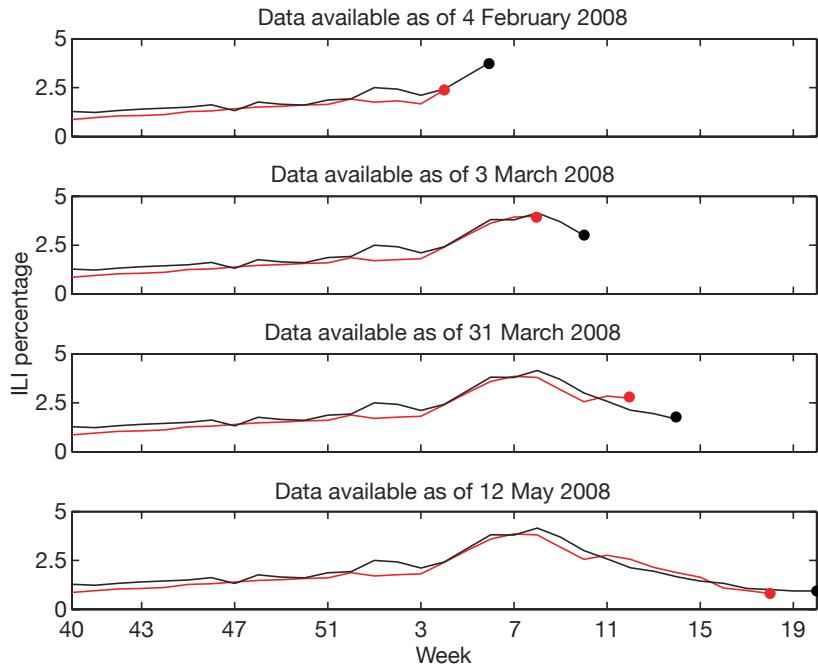
Platform effects

Hierarchy of limitations in machine learning

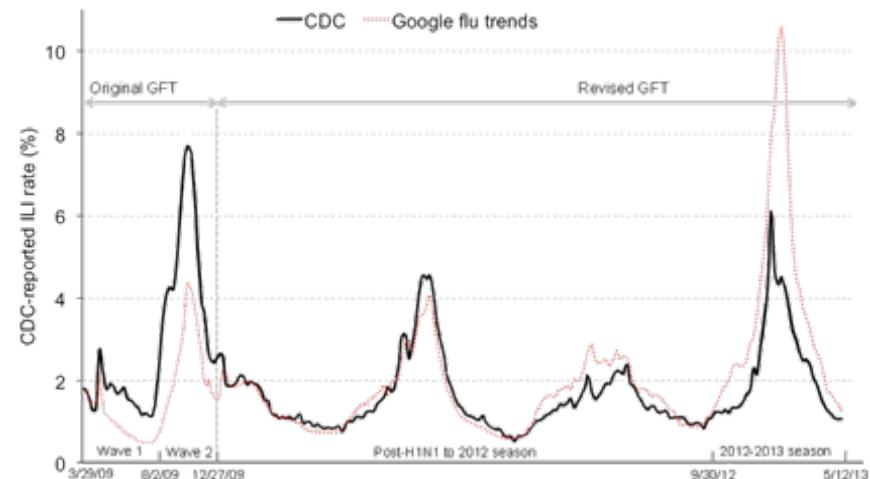
Problems of cross-validation

Summary and conclusion

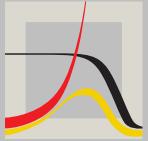
References



Ginsberg et al., 2012, *Nature*

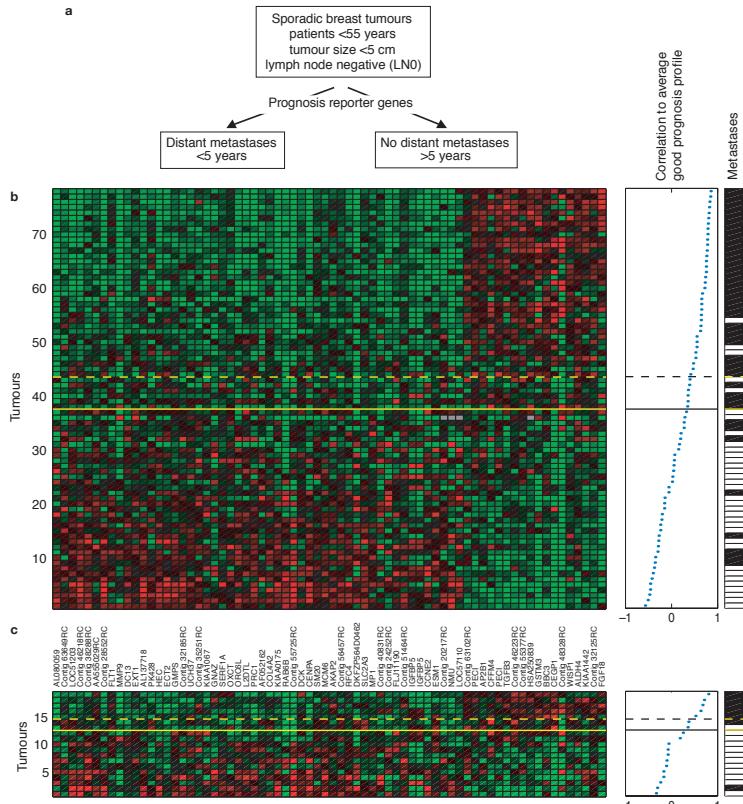


Santillana et al., 2014, *Am. J. Prev. Med.*



Real-world testing of ML results

- van't Veer et al. (2002) found 70 genes correlated with developing breast cancer
 - Of course the correlations were optimal, post-hoc. But did it generalize?





Implementation testing

Introduction

Brief historical tour

Bias in geotagged tweets

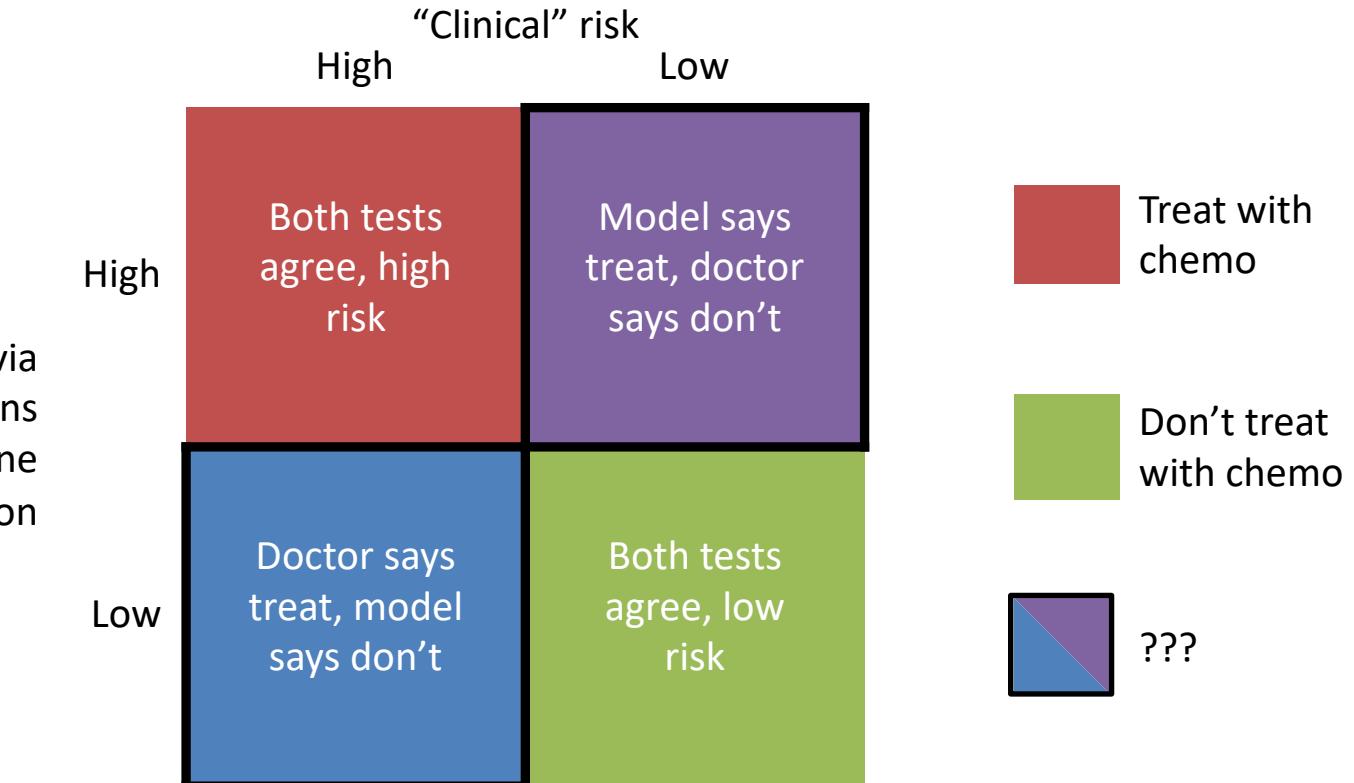
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Cardoso et al., 2016, NEJM



Implementation testing

Introduction

Brief historical tour

Bias in geotagged tweets

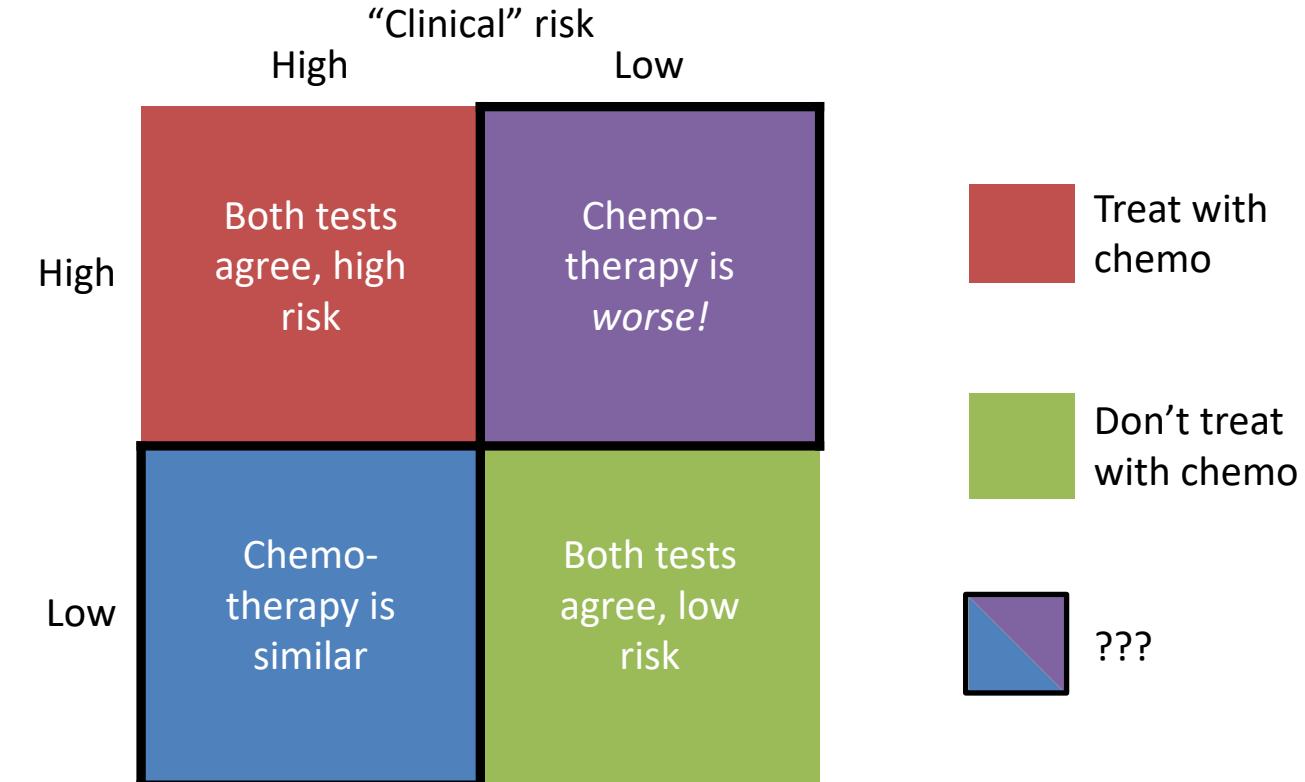
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Cardoso et al., 2016, NEJM



Implementation testing

Introduction

Brief historical tour

Bias in geotagged tweets

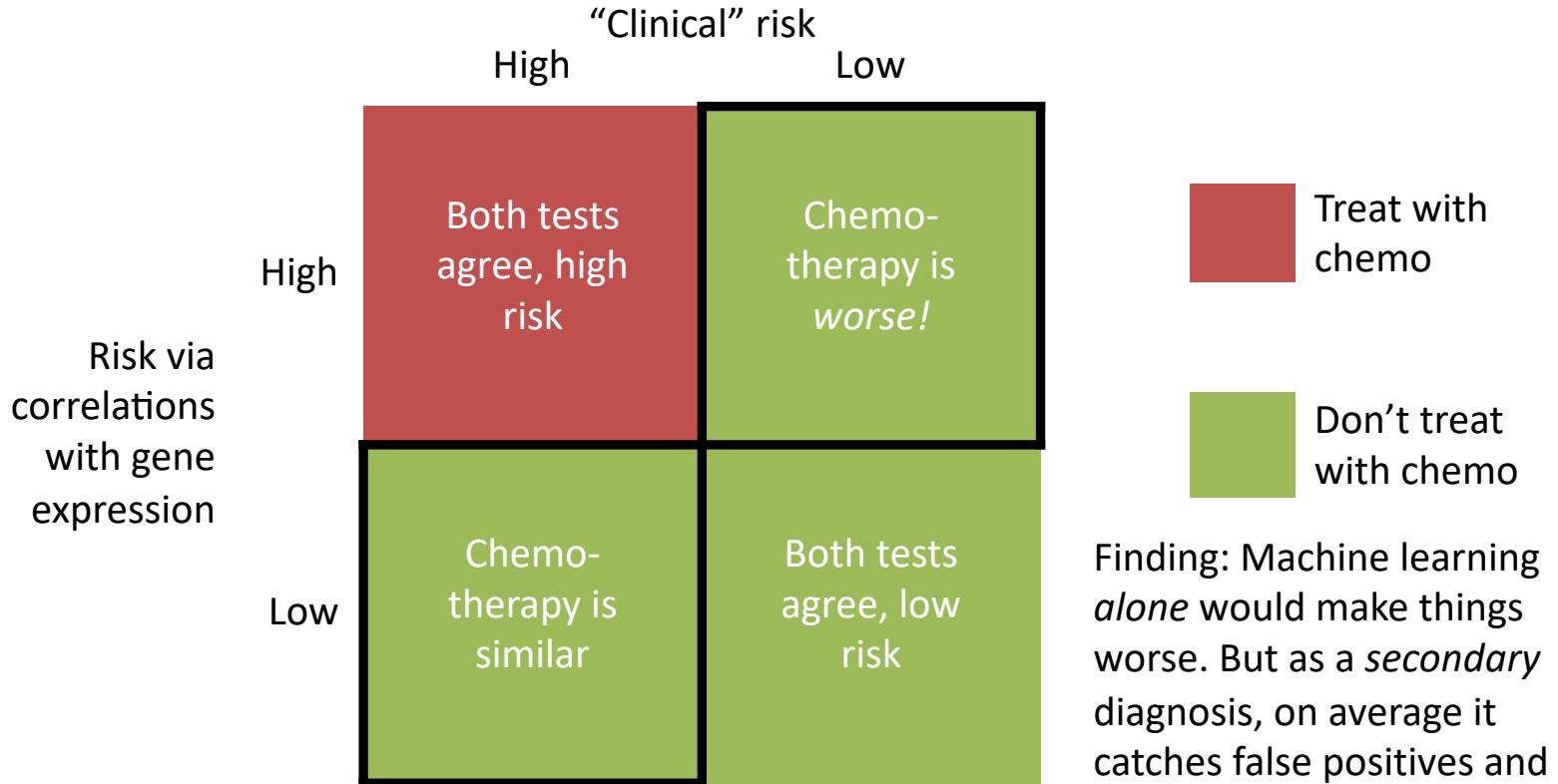
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Cardoso et al., 2016, NEJM



Implementation testing: Details

Introduction

Brief historical tour

Bias in geotagged tweets

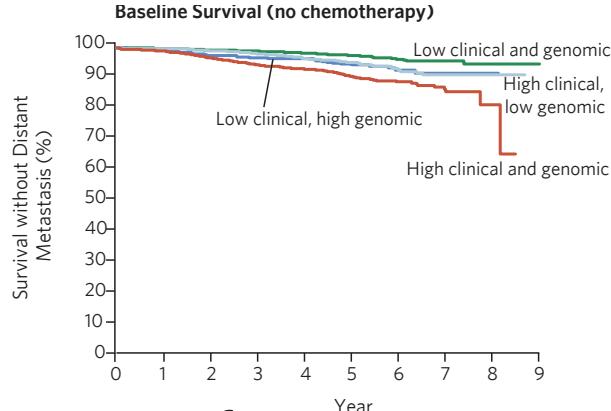
Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

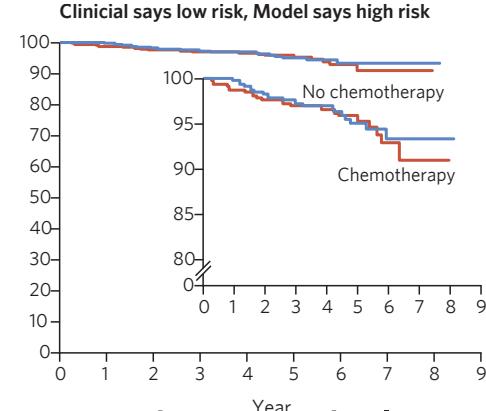
Summary and conclusion

References

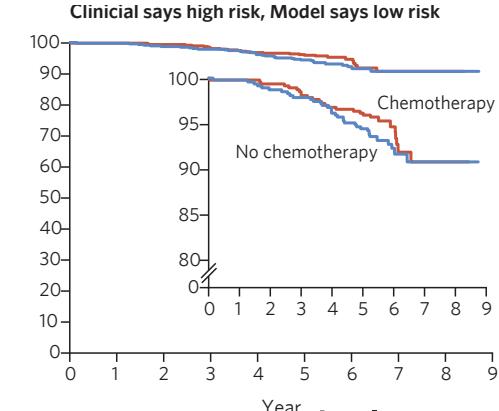


- Before experiment (training data)

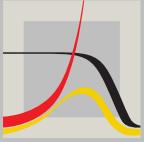
(Note: still limitations in how experimental subjects may be unrepresentative.)



- High model risk, low clinical risk: randomize. Chemo worse!



- Low model risk, high clinical risk: chemo makes no difference



Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

Summary and conclusion



Summary

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

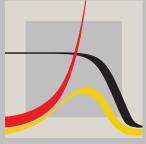
Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

- Biases exist, but are not simple, and may be unknowable in general
 - When we have a comparison source, we can calibrate, but better may be to find appropriate use cases
- Commercial platforms are not always fit for research; but we can try to investigate how their design and incentives
- Machine learning presents new opportunities, but has multiple failure points (often corresponding to long-recognized problems) that must be recognized and dealt with
 - Prediction, and cross validation, have fragilities
 - Out-of-sample testing is always a good idea



Conclusion

- There are problems with data; these have been widely recognized, and we are making progress on how to work with new forms of data, including opaque secondary and commercial data
- There are still fundamental problems limitations of different modeling approaches, and how modeling relates to the world: machine learning is supercharging these, forgetting lessons of the past

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References



Other work of mine

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References

- This is what I think is most interesting to this audience
- I have other work:
 - Trying to use modeling to *imagine* alternative states of the world (Richardson, Malik et al., 2021)
 - Why do “technical” people have such a narrow view of the world, and how do some come to change? (Malik & Malik, 2021)
- My current work is on “AI ethics”, which I take to mean, how do we rigorously and responsibly develop and deploy (or choose not to develop or deploy) modeling, applied to large-scale data? How do we choose what modeling approach is appropriate (if any)?
 - This is largely aimed at biomedical use cases, but developed guidance will, I hope, apply to any field of social science



References

Introduction

Brief historical tour

Bias in geotagged tweets

Platform effects

Hierarchy of limitations in machine learning

Problems of cross-validation

Summary and conclusion

References