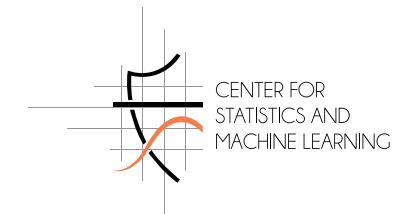


When (and why) we shouldn't expect reproducibility in machine learning-based science: Culture, study design, causality, and asymptotics

Momin M. Malik

Senior Data Science Analyst - AI Ethics, Center for Digital Health, Mayo Clinic
School of Social Policy & Practice, University of Pennsylvania
Institute in Critical Quantitative, Computational, & Mixed Methodologies

Presented at: The Reproducibility Crisis in ML-Based Science
Center for Machine Learning and Statistics, Princeton University
2022 July 28 [online]





Outline

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- Methodological problems:
 - Sampling frame and measurement: ML tries to circumvent internal/construct validity, and sampling frame. Sometimes this works. Sometimes it doesn't.
 - "Prediction" vs. causality: "prediction" has a lot of complexity that casual usage ignores
 - Model metrics are estimators! They have asymptotic distributions, can be biased, etc.
 - Dependencies are a form of leakage, and bias the CV estimators of model success
 - We should figure out asymptotic distributions to help design tests and power calculations, and start using them
- Contextual points:
 - Cultural issues: Lack of exposure to the entirety of research methods. To a certain extent, expecting replication might be understanding science in a bad way
 - Lessons from other fields: We probably won't get reform until we have a crisis, and we probably won't get to a crisis



Background (Malik, 2020)

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

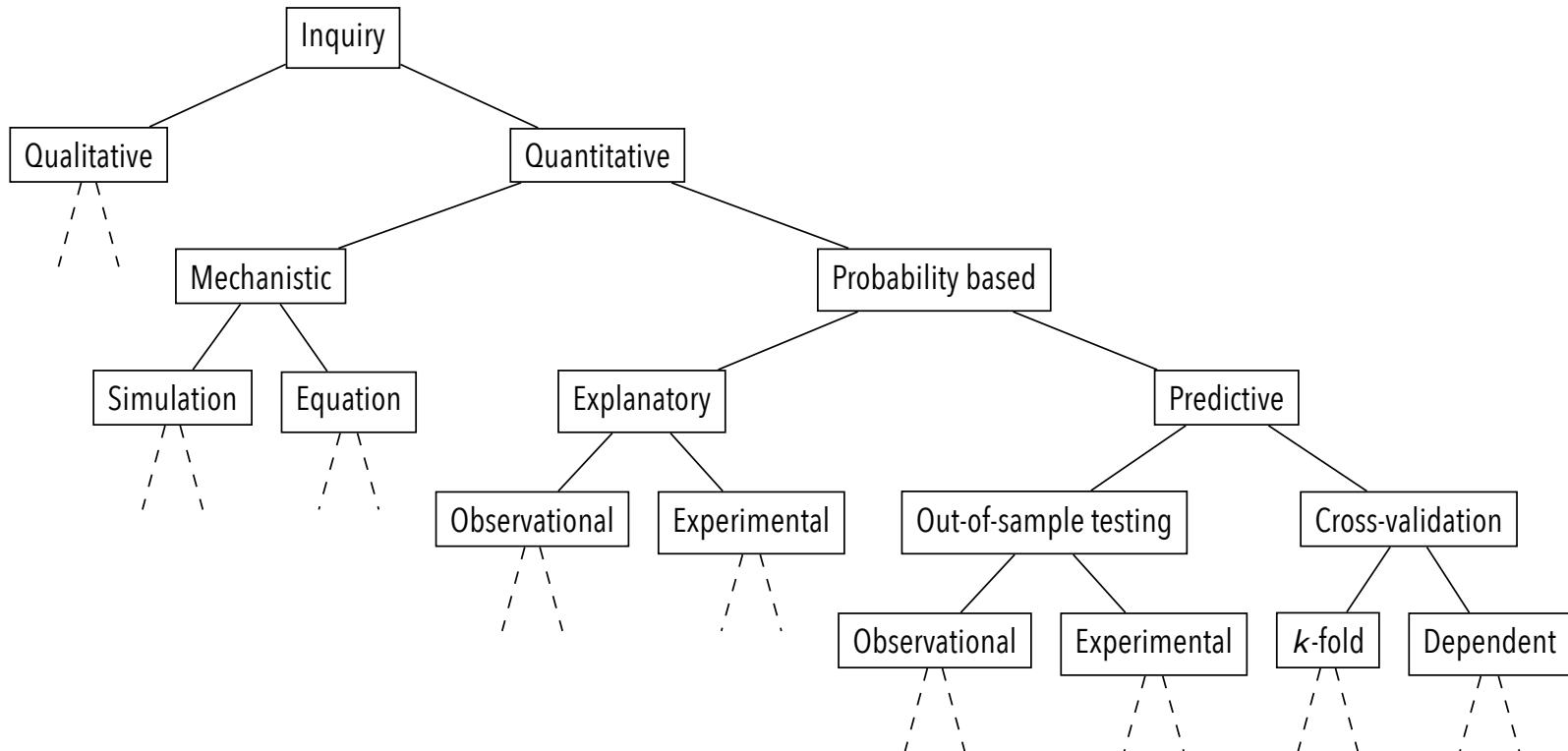
Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code





Introduction

**Sampling
frame and
measurement**

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

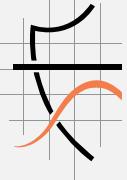
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

Problem 1: Sampling frame and measurement



Traversing the hierarchy of limitations

Introduction

**Sampling
frame and
measurement**

"Prediction"
vs. causality

Model metrics
as estimators

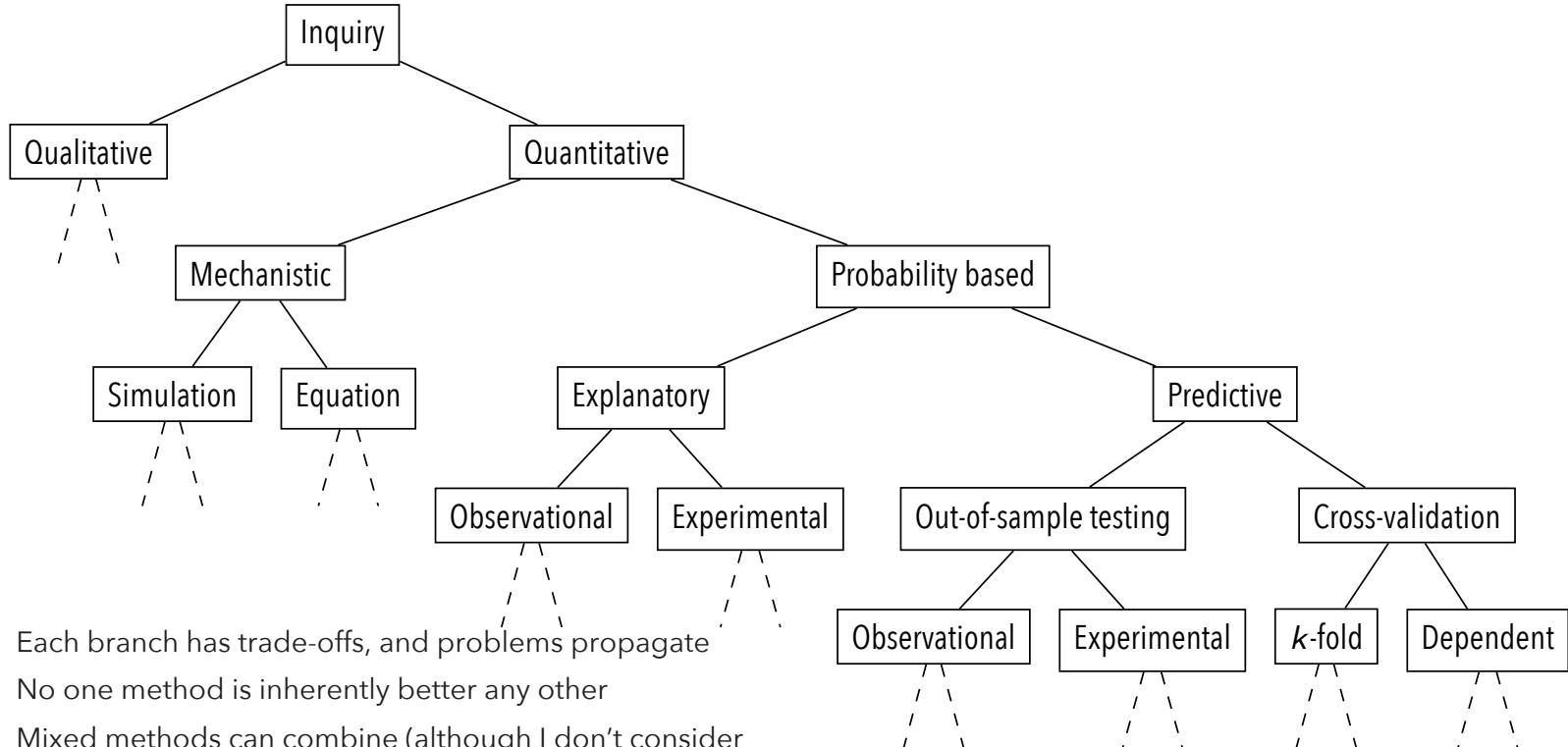
Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code



- Each branch has trade-offs, and problems propagate
- No one method is inherently better than any other
- Mixed methods can combine (although I don't consider this in the paper)



Quantification locks in meaning

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

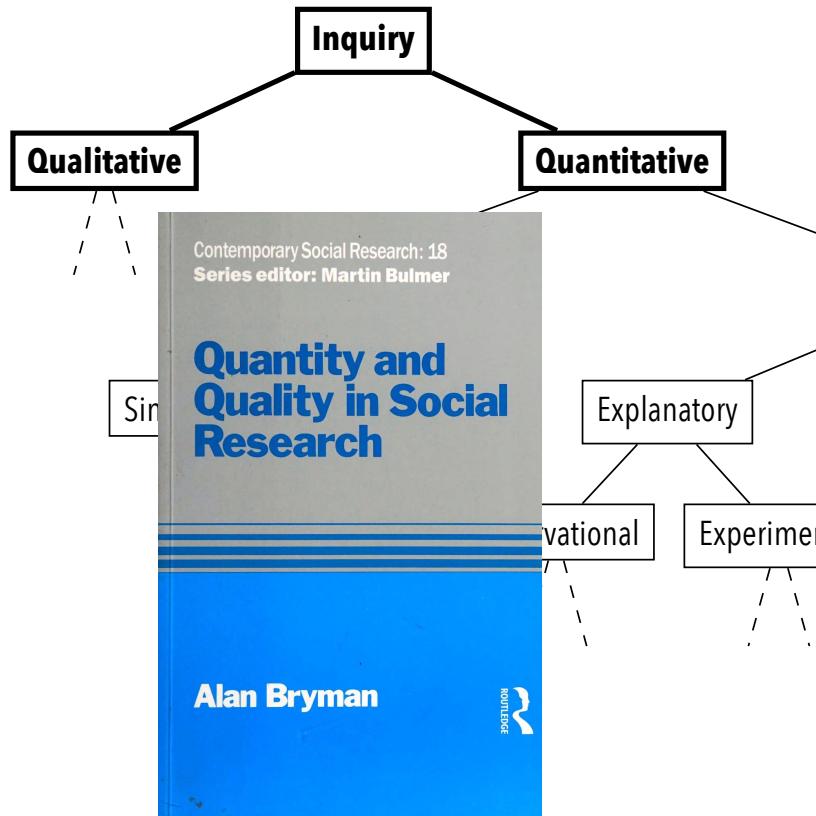
Cultural issues

Lessons from other fields

Summary and conclusion

References

Appendix:
Simulation code



- 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

**Sampling
frame and
measurement**

"Prediction"
vs. causality

Model metrics
as estimators

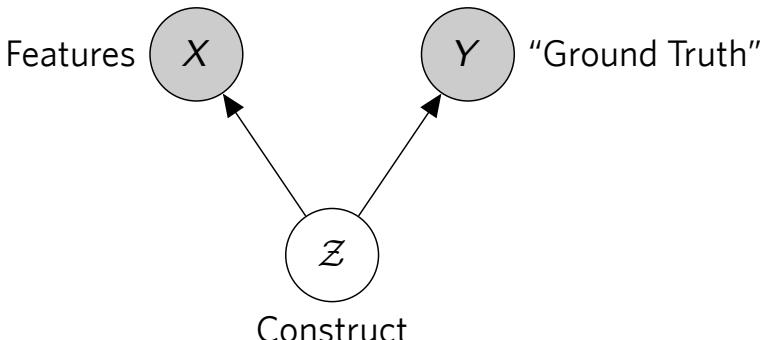
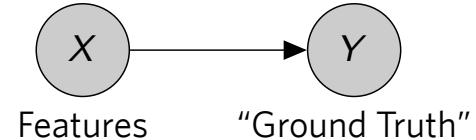
Cultural issues

Lessons from
other fields

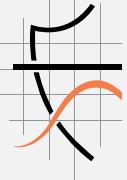
Summary and
conclusion

References

Appendix:
Simulation
code



- *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), and even can be gamed



Example: Epic sepsis model

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

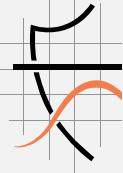
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- Wong et al. (2021) found that a model to predict sepsis from the electronic health records company Epic worked far less well than claimed
- AUC of .63, versus what Epic reported of .76 to .83
- One possible culprit: *different definitions*. Epic developed its model based on defining sepsis by the point where physicians intervened (what there was direct data for). Wong et al.'s evaluation was based on defining sepsis by meeting a certain number of CDC and ICD-10 criteria
- Of course the model as fitted wouldn't generalize! Maybe the same model, re-fitted on the "better" measure, would work



Stats and ML use central tendencies

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

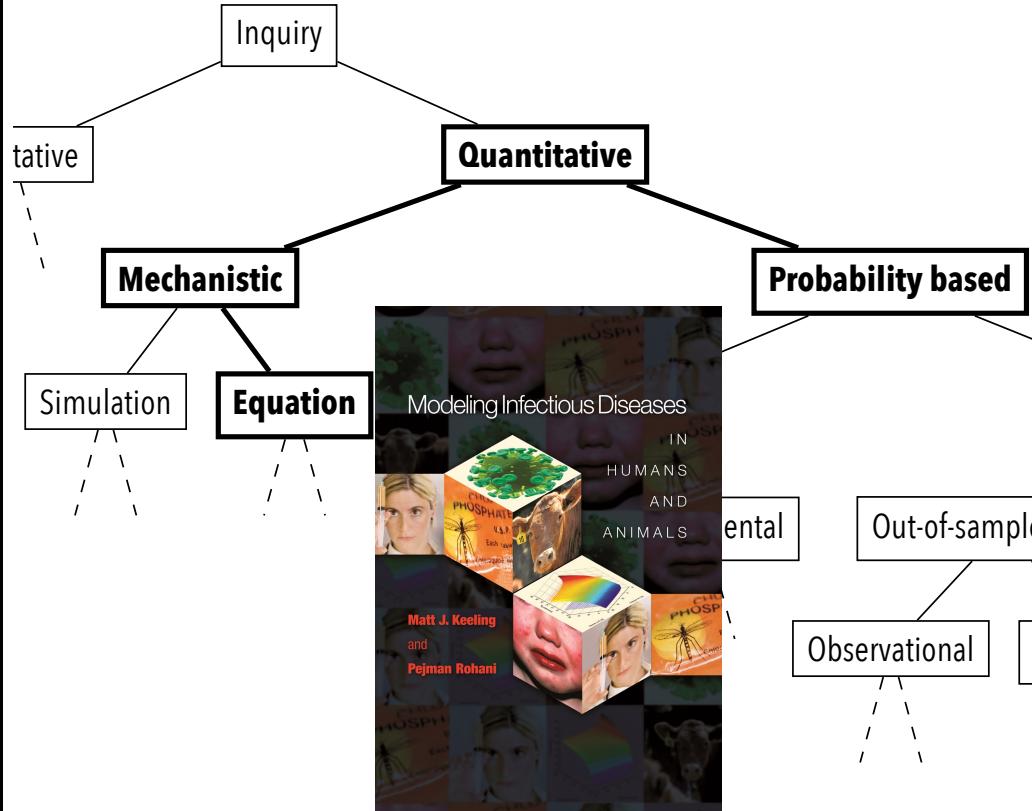
Cultural issues

Lessons from other fields

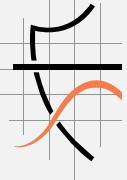
Summary and conclusion

References

Appendix:
Simulation code



- 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 (we cannot do anything with an n of 1!)



Importance of sampling frame

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- Because ML uses the same fundamental mechanism as stats (reducing aggregates via central tendency), it has the same issue that *results will only generalize insofar as the sample is representative*
 - (That was the problem with "Dewey defeats Truman", and also the "Literary Digest poll" in survey sampling, that led to reform there)
- The "patterns" we "recognize" are correlations, not necessarily universal regularity, so we can't ignore the sampling frame
- "Sampling on the dependent variable" is a classic problem: Cohen and Raths (2013) have an amazing *mea culpa* where they note that they filtered Twitter users to only those who had a signal for political orientation. That was an unrealistic sampling frame



Core issues of study design

Introduction

**Sampling
frame and
measurement**

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

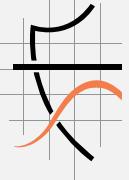
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- Sampling frame is typically taught in social sciences, not necessarily in machine learning
- Measurement models are the domain of psychometrics, and are almost completely unknown in ML (Jacobs & Wallach, 2019)
- These are a standard part of education that ML should make room for (will return to later under "culture")



Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

Problem 2: "Prediction" vs. causality



Causality is hard, maybe too hard

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

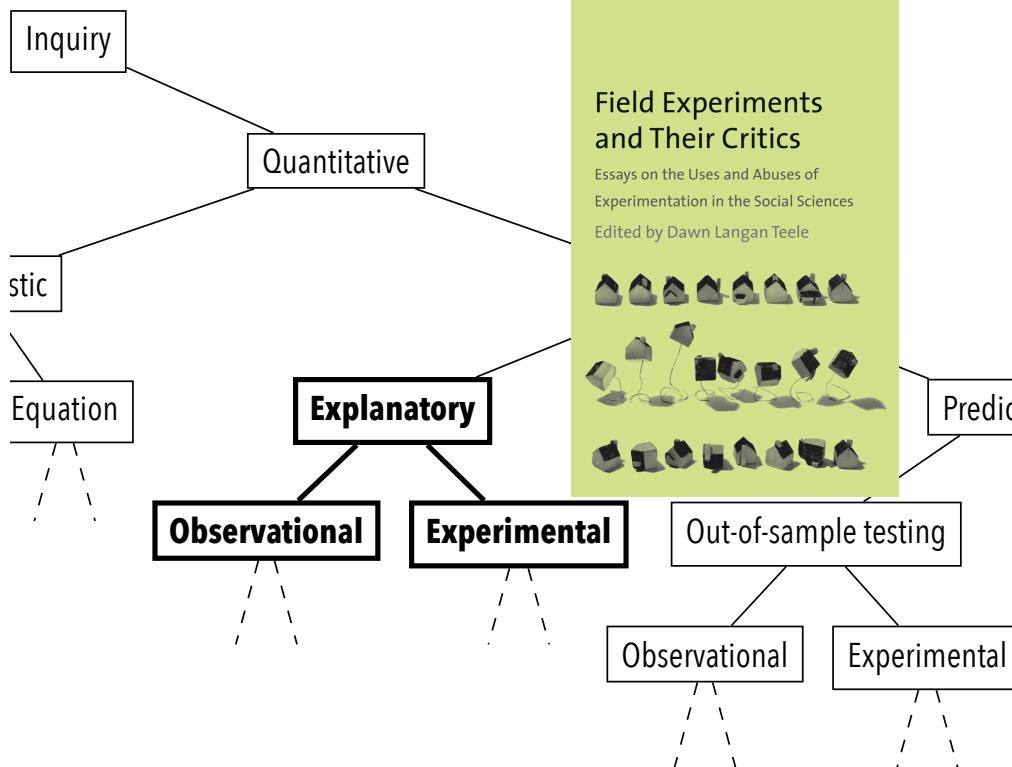
Cultural issues

Lessons from other fields

Summary and conclusion

References

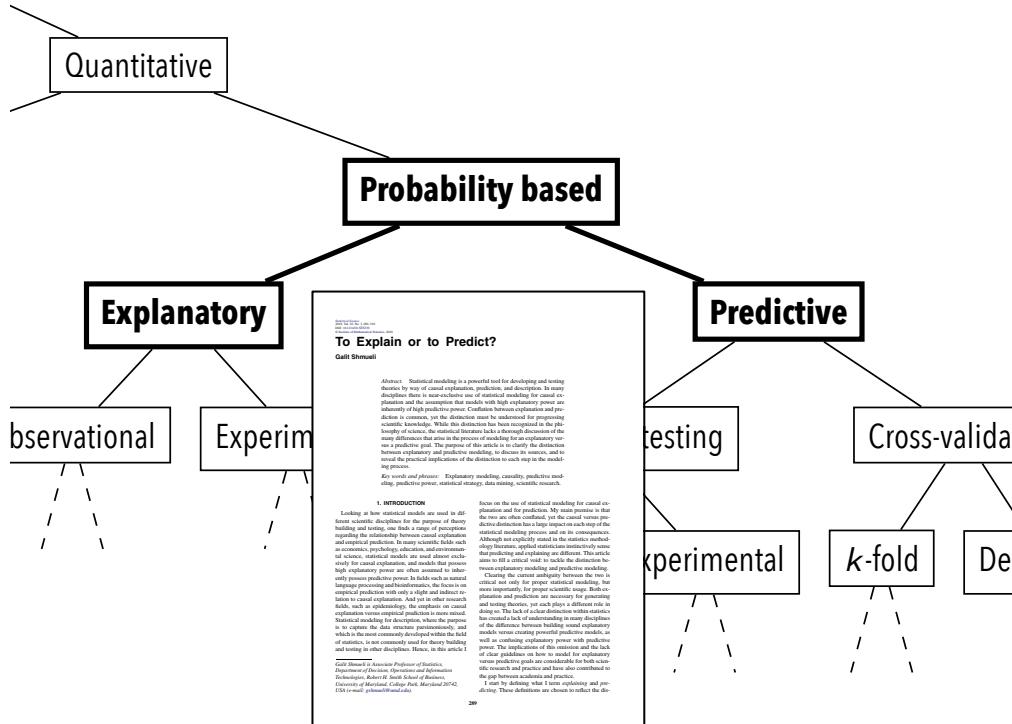
Appendix:
Simulation code



- 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, 2010)



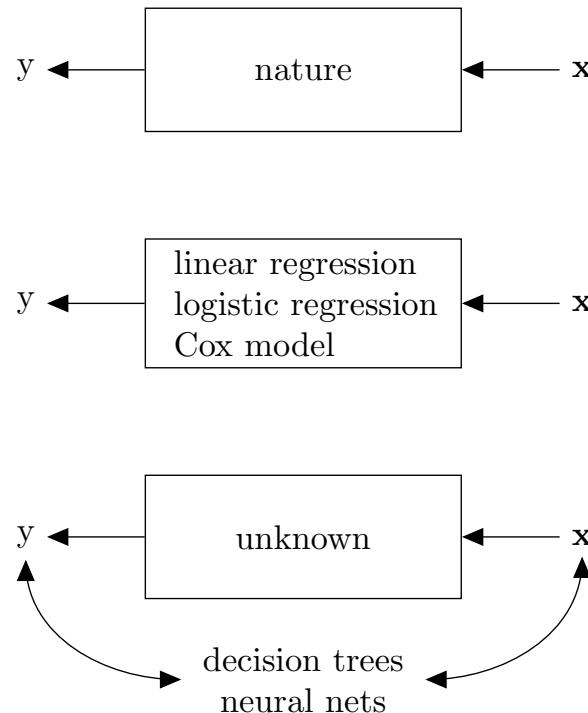
ML is “prediction” only



- “Predictions” are defined as what minimizes loss *within a predetermined frame*
 - *Correlations* do this
 - Non-causal correlations can sometimes predict well within a frame, but they frequently don’t explain, and can fail outside
 - If that was the definition (Milton Friedman: “prediction in the presence of change”), correlations wouldn’t work, but that is hard to formalize



A “realist” definition for machine learning



- Realist definitions: what things are, rather than what they aspire to be
- 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)



Leads to two separate goals

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

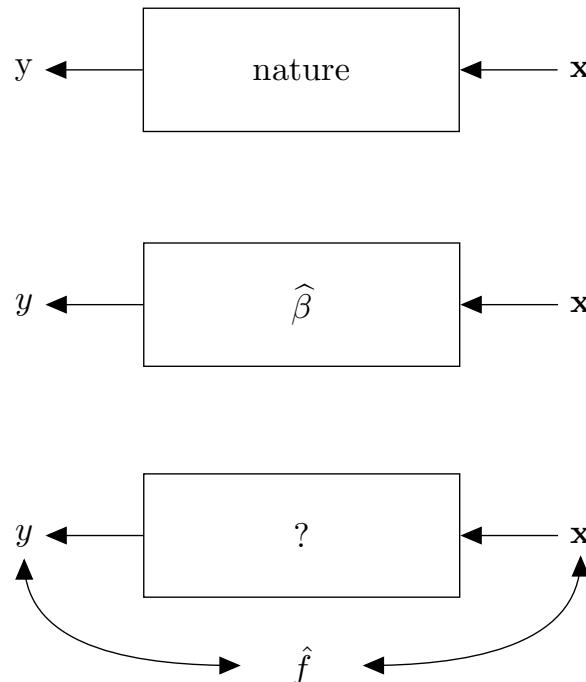
Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

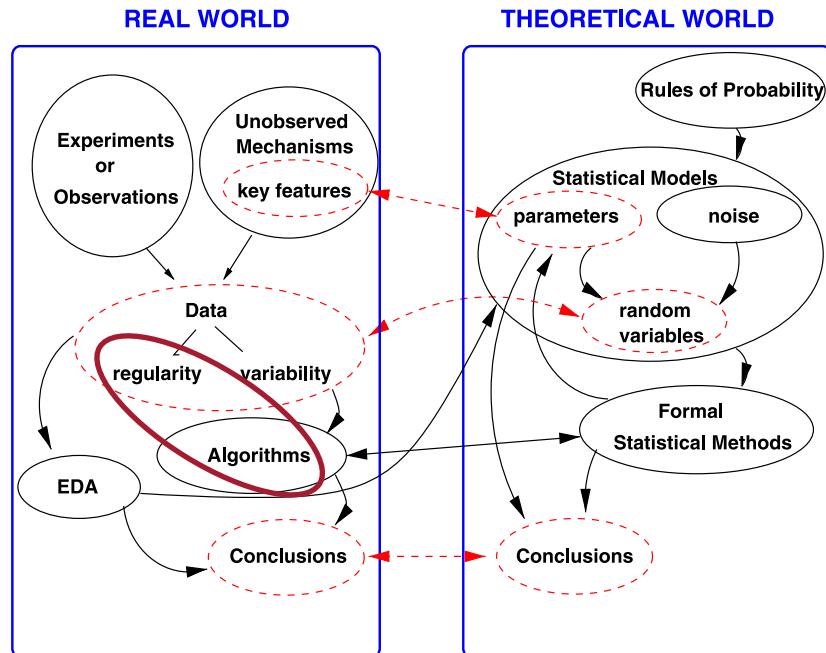


- Non-causal ("spurious") correlations may fit robustly (e.g., latent common cause)
 - Breiman, 2001: "prediction problems"
 - Shmueli, 2010: "to predict"
 - Kleinberg et al., 2015: "umbrella problems"
 - Mullainathan & Spiess 2017: "**y-hat problems**"
- Carefully built models that capture causality (or "pure" associations) may fit poorly overall
 - Breiman: "information"
 - Shmueli: "to explain"
 - Kleinberg et al.: "rain dance problems"
 - Mullainathan & Spiess: "**beta-hat problems**"



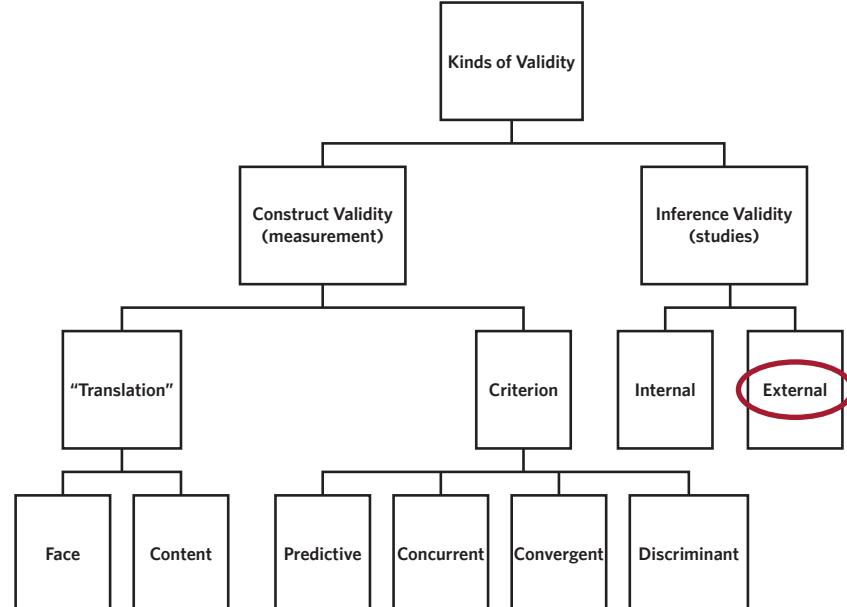
ML: Only external validity

- Introduction
- Sampling frame and measurement
- "Prediction" vs. causality
- Model metrics as estimators
- Cultural issues
- Lessons from other fields
- Summary and conclusion
- References
- Appendix:
Simulation code



Kass, 2011

Adapted from Borgatti, 2012





Levels of prediction (Rescher, 1998)

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

Cultural issues

Lessons from other fields

Summary and conclusion

References

Appendix:
Simulation code

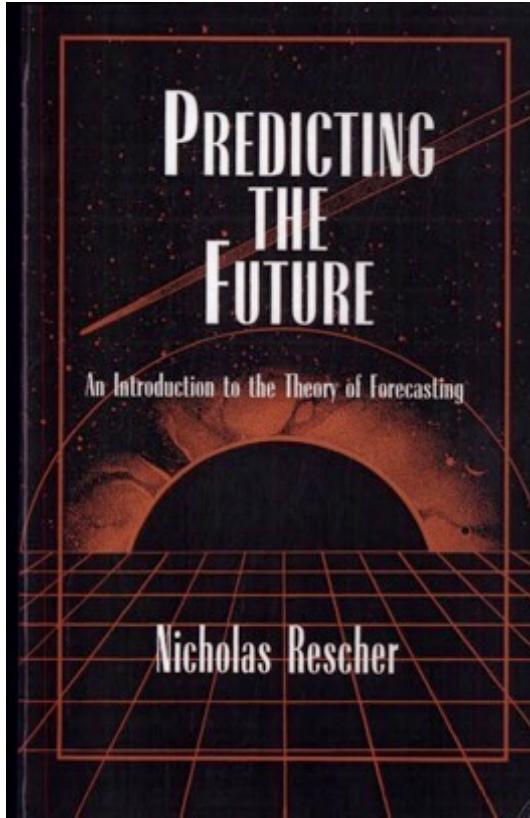


TABLE 6.1: A SURVEY OF PREDICTIVE APPROACHES

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



"Things do change" (Hoadley, 2001)

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

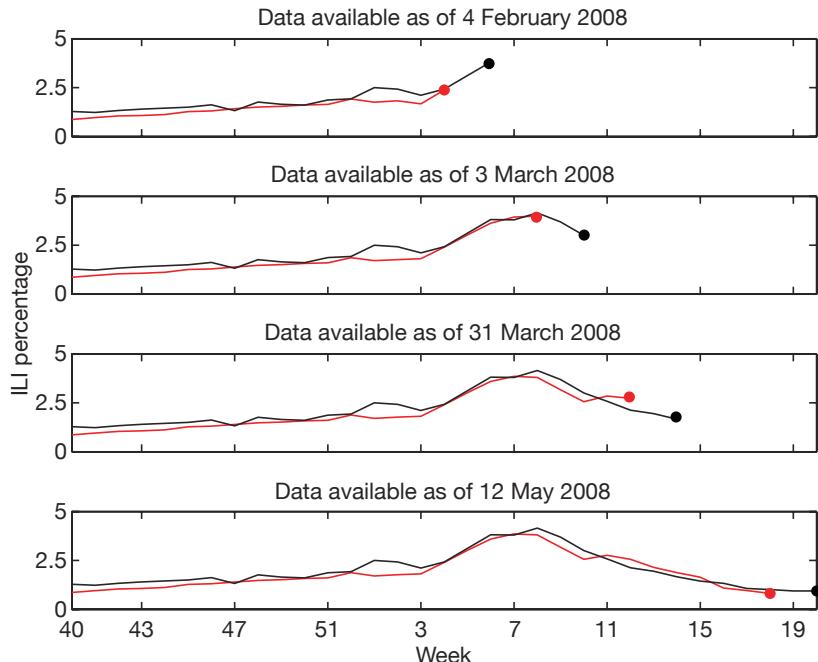
Cultural issues

Lessons from other fields

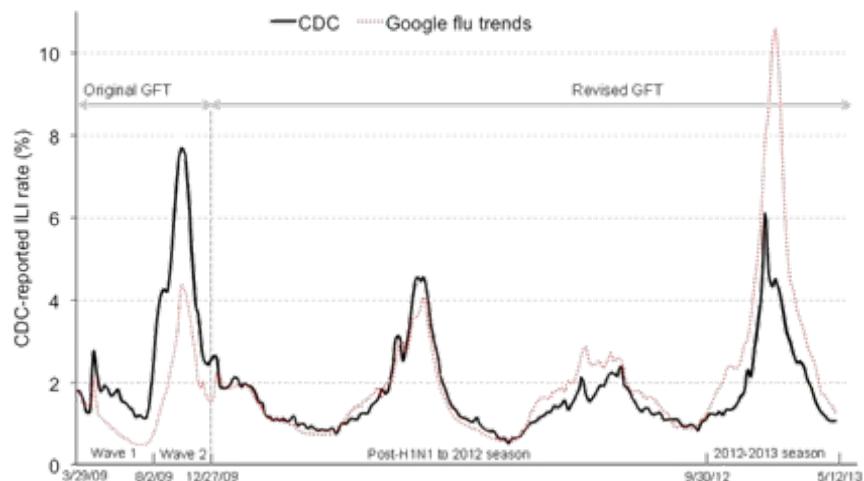
Summary and conclusion

References

Appendix:
Simulation
code



Ginsberg et al., 2012



Santillana et al., 2014



Correlations can't "predict in the presence of change" or interventions

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

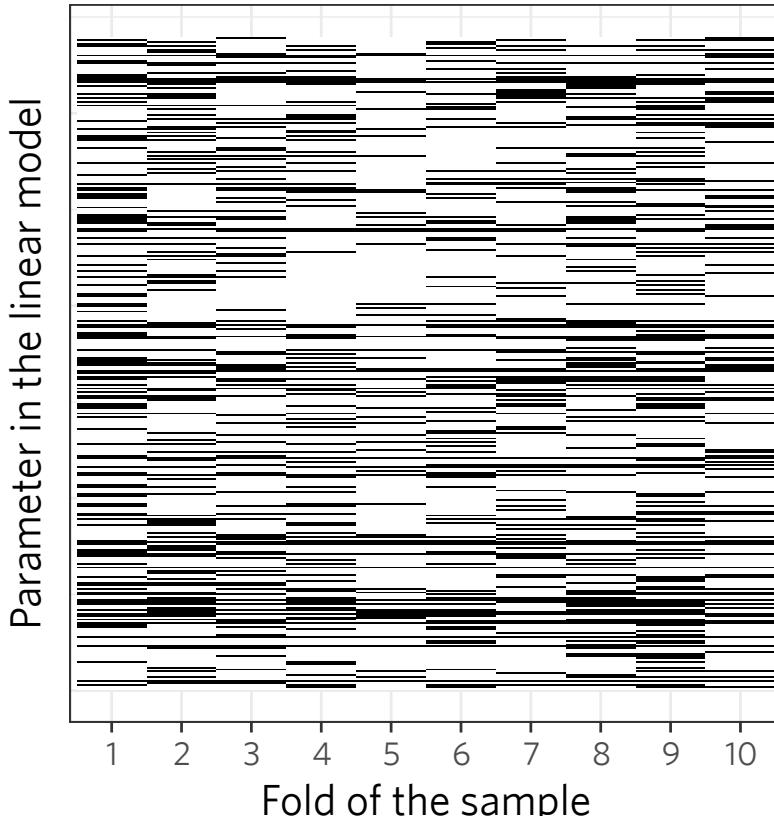
Cultural issues

Lessons from other fields

Summary and conclusion

References

Appendix:
Simulation code



- Very different sets of correlations can "predict" (correlate) equally well (Mullainathan and Spiess 2017)
 - Breiman (2001) called this the "Rashomon Effect"
- But different fits suggest very different outputs under covariate shift, and under interventions



Pet peeve: language

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

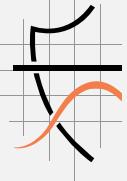
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- Communication: **stop saying "prediction" if it is really "correlation"**
 - **The use of 'prediction' leads to false, inflated expectations.** Instead of saying prediction for post-hoc demonstrations (Gayo-Avello, 2012), use "retrodiction" instead: it is awkward, but that's what we need. For time series: nowcasting, back-testing.
 - Attempts to model partial correlation, i.e., "ceteris paribus", can be described with "association"
- Prediction is overused as it is
 - Statements like "predict the probability of risk", or "calculate the probability of a likelihood" are redundant if not nonsensical (akin to, "a probability of a probability [of a probability]").
 - Probabilities and risks are already latent, predictions are of things that will manifest. We should say that *estimate* probabilities and risk (say *estimated probabilities*, etc.)
 - Use "detection" or "classification" if labels are manifest but unknown. E.g., we don't "predict" race. "Detecting" and "predicting" cancer imply two very different tasks.



Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

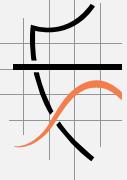
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

Problem 3: Model metrics are estimators, with unknown asymptotics and sources of bias



Model metrics as estimators

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

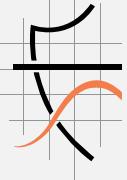
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- If we make a commitment to a statistical view of the world (unobservable but inferable underlying regularity realized with haphazard variability), then the precision, recall, AUC, etc., are *estimators* of the underlying quantity of out-of-sample performance
 - Quantifying uncertainty provides a hedge on performance claims
- We can frame and study their properties statistically!
 - We can design corrections, find instrumental variables, etc.
 - *Dependencies* cause test error to be biased (and, in a simple case, error has a generalized non-central chi-square distribution, which is heavily right-tailed, versus the symmetry of a binomial distribution)
 - Metrics other than accuracy (binomial) look like they have weird distributions. Somebody should look into this, and also design asymptotic tests
 - Under this view, it makes sense to use instrumental variables to try and get unbiased estimates of out-of-sample performance! (Kleinberg et al., 2018)



Matrix bias-variance decomposition

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

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

irreducible ("Bayes") error bias squared variance of the estimator "optimism"



Classic argument for CV

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

Training:

$$\begin{aligned}\text{err}(\hat{\mu}) &= \frac{1}{n} \mathbb{E}_f \|Y - \hat{Y}\|_2^2 \\ &= \frac{1}{n} \left[\text{tr} \Sigma + \|\mu - \mathbb{E}(\hat{Y})\|_2^2 + \text{tr} \text{Var}_f(\hat{Y}) - 2 \text{tr} \text{Cov}_f(Y, \hat{Y}) \right]\end{aligned}$$

Testing:

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

The difference is the *optimism* (Efron, 2004; Rosset & Tibshirani, 2018):

$$\text{Opt}(\hat{\mu}) = \text{Err}(\hat{\mu}) - \text{err}(\hat{\mu}) = \frac{2}{n} \text{tr} \text{Cov}_f(Y, \hat{Y})$$



Apply this to non-iid data

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- 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!

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



One draw as an example

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

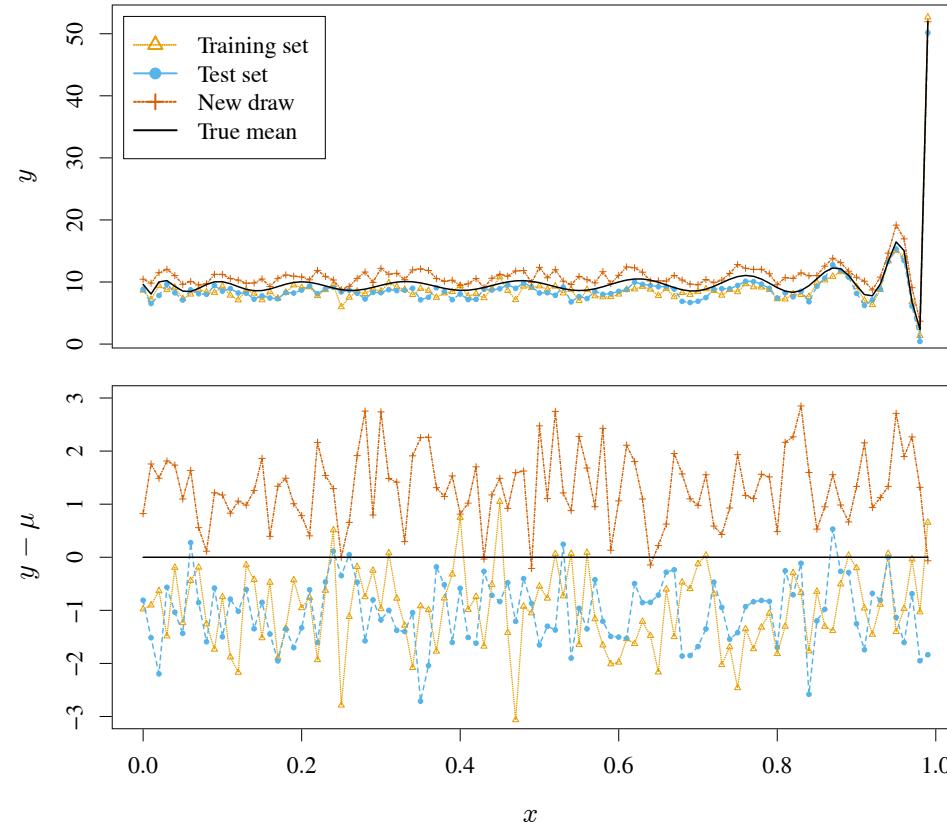
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

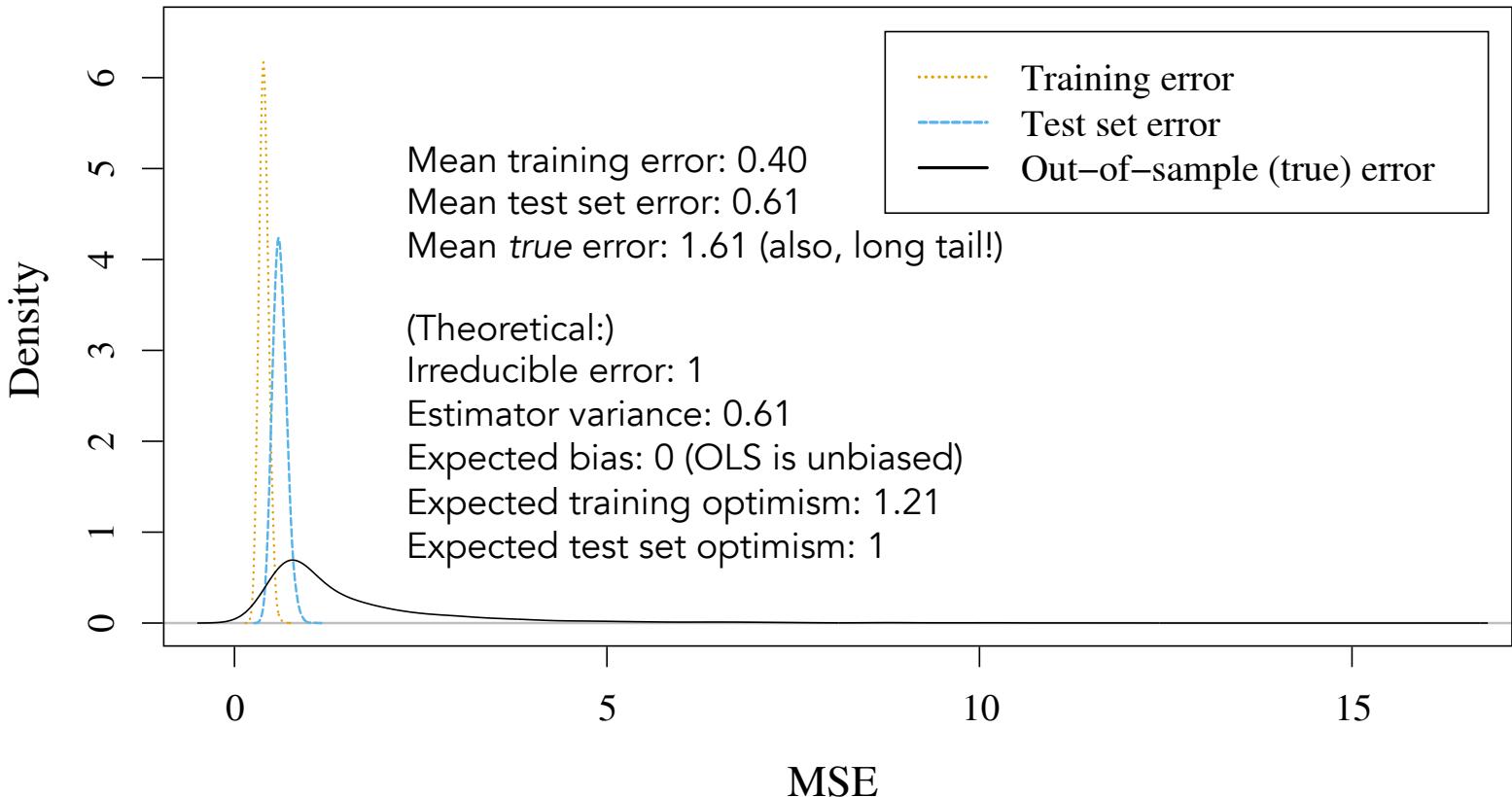
Correlation
between
observations can
pull training and
test
observations
close to one
another, but
potentially far
from an
independent
draw





Simulated MSE

- Introduction
- Sampling frame and measurement
- "Prediction" vs. causality
- Model metrics as estimators**
- Cultural issues
- Lessons from other fields
- Summary and conclusion
- References
- Appendix:
Simulation code





Lessons: Split by dependencies

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- ML needs to contend with dependencies, because the iid assumption matters for estimates of model performance
 - Even statistical relational learning, which models dependencies, doesn't talk about data splitting by dependencies
- Maybe we can't make a better *model*, but dependencies are a form of leakage between training and test sets
 - We can use the framework of "optimism" to understand and quantify this (meta-meta-prediction is useful; Rescher, 1998)
 - Ideally, no dependencies between training and test sets
 - Unfortunately, the mean function and covariance function are jointly unidentifiable nonparametrically (Opsomer & Jean, 2011), so we will have to rely on theory and limited explorations (e.g., ACF, PACF)



How are metrics distributed? (Preliminary explorations)

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

Cultural issues

Lessons from other fields

Summary and conclusion

References

Appendix:
Simulation code

- Under this specification and DGP, the test error has a "generalized non-central chi-squared" distribution
- But even in the iid case, we know frighteningly little about asymptotic distributions (in that I found no work other than around accuracy, which is binomial) and the variability they might suggest
- A quick simulation of a logistic fit to $X_i \sim \mathcal{N}(0, 1)$ and $Y_i \sim \text{Bin}(\text{logistic}(x_i))$ at $n = 101$ (small, prime-numbered sample size) and $n = 10,000$ (large sample size) gives reasons for worry



Distributions of counts? $n = 101$

($n_{\text{sim}} = 50,000$)

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

Lessons from
other fields

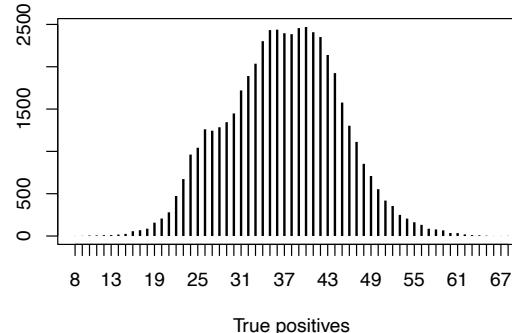
Summary and
conclusion

References

Appendix:
Simulation
code

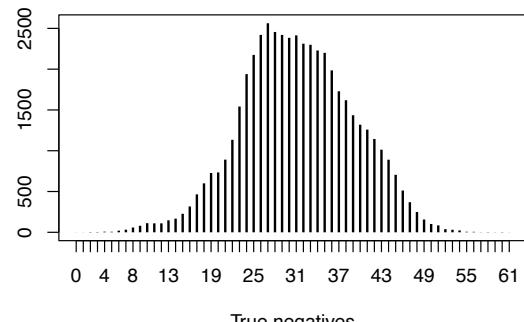
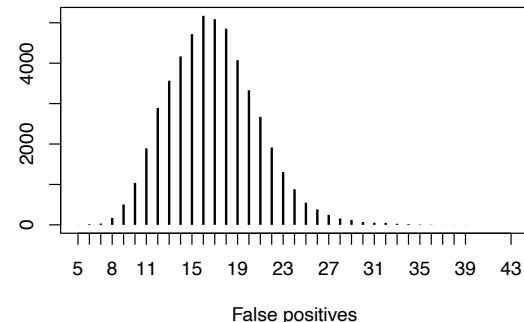
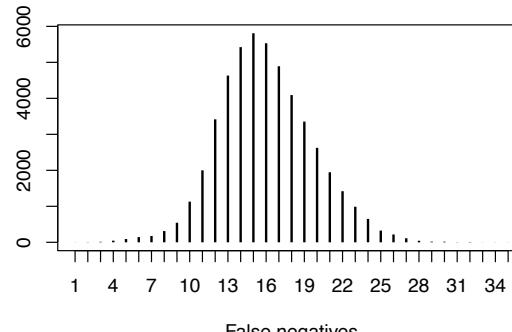
Actual positive

Predicted positive



Actual negative

Predicted negative





Distribution of precision/recall? $n = 101$ $(n_{\text{sim}} = 50,000)$

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

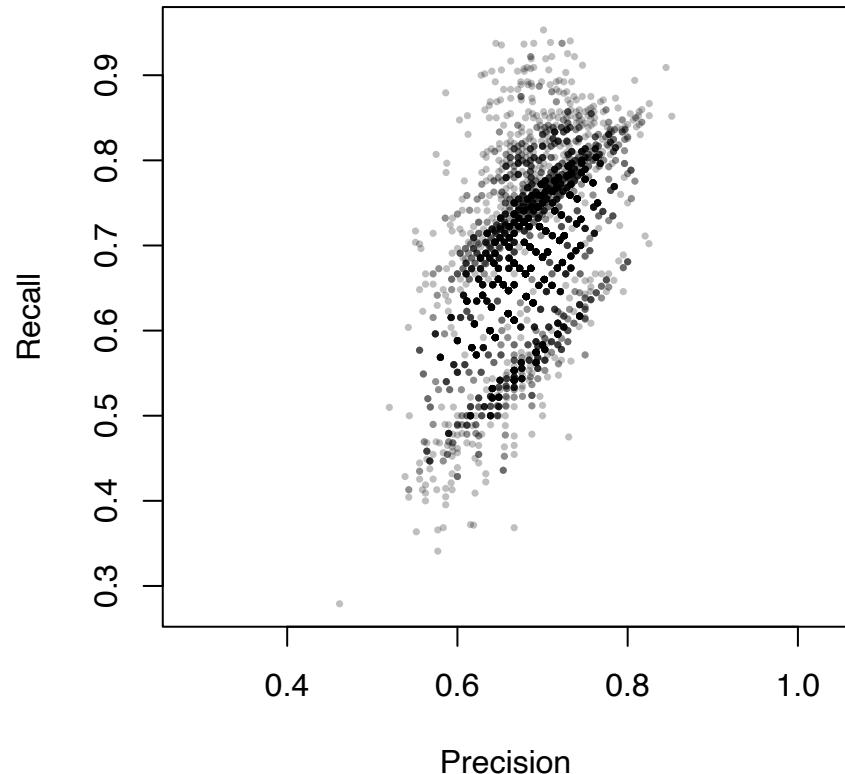
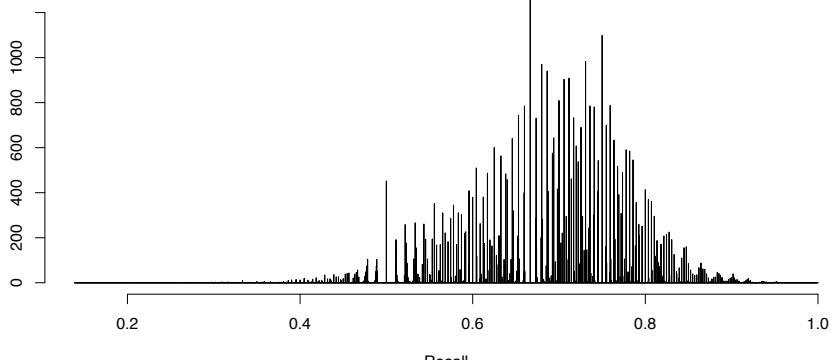
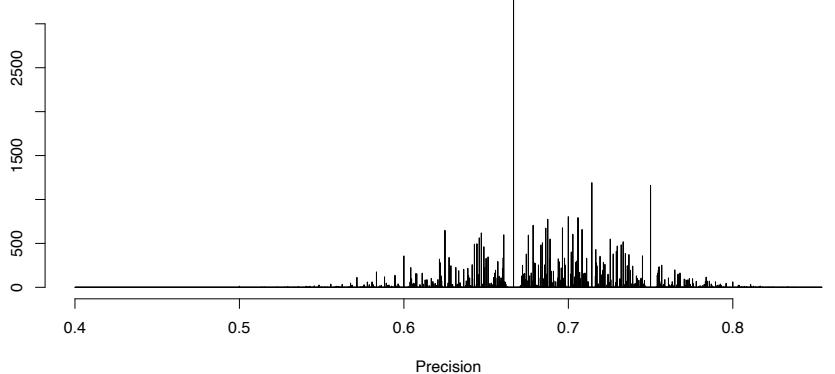
Cultural issues

Lessons from
other fields

Summary and
conclusion

References

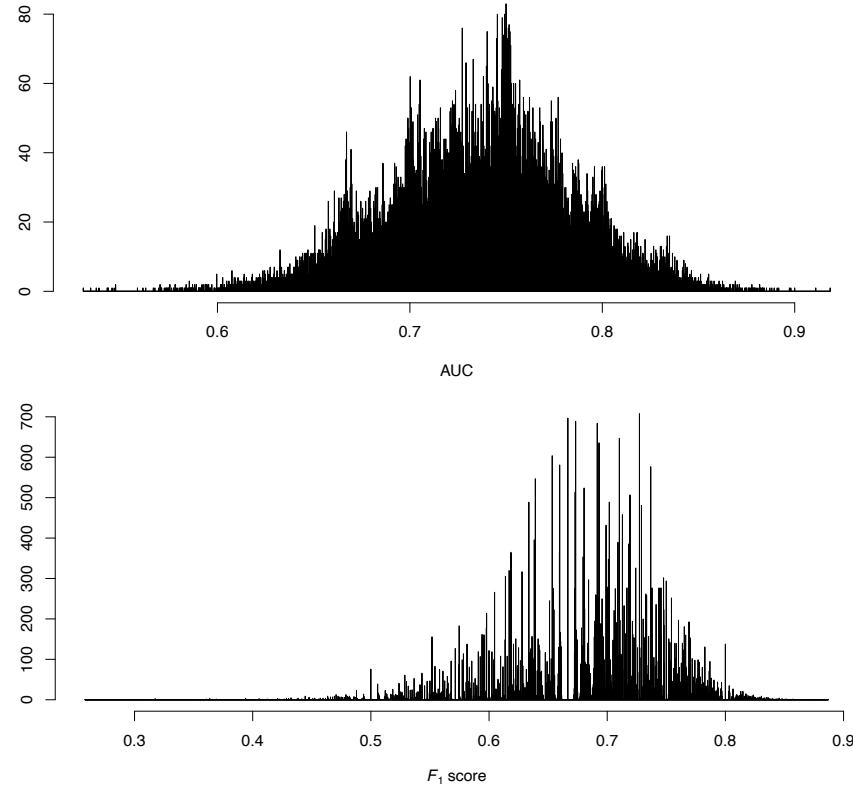
Appendix:
Simulation
code





Distribution of AUC/ F_1 ? $n = 101$ ($n_{\text{sim}} = 50,000$)

- AUC has a 95% empirical confidence (tolerance) interval of [.64, .83]
- Small sample size is an issue, but also, we usually try to stay away from modeling ratios, since the ratio of two normal distributions is Cauchy, which has no expected value (practically: ratios are unstable)
- What about larger sample size?





Distributions of counts? $n = 10^4$

($n_{\text{sim}} = 50,000$)

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

**Model metrics
as estimators**

Cultural issues

Lessons from
other fields

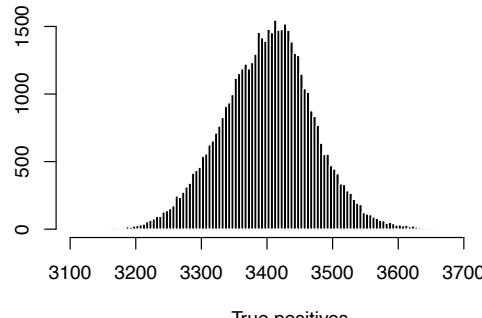
Summary and
conclusion

References

Appendix:
Simulation
code

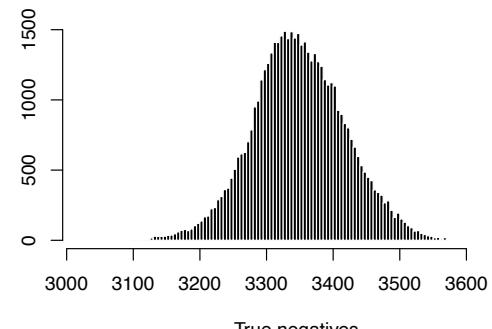
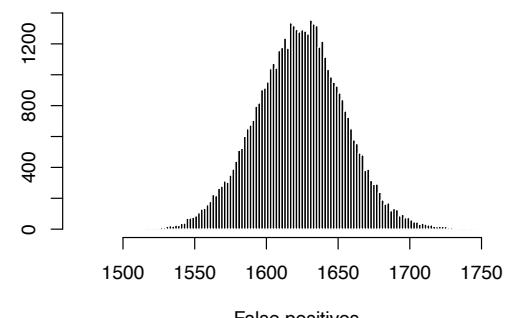
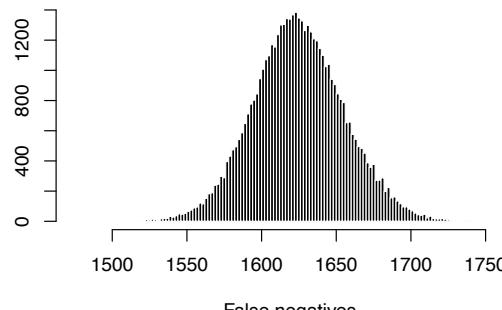
Actual positive

Predicted positive



Actual negative

Predicted negative

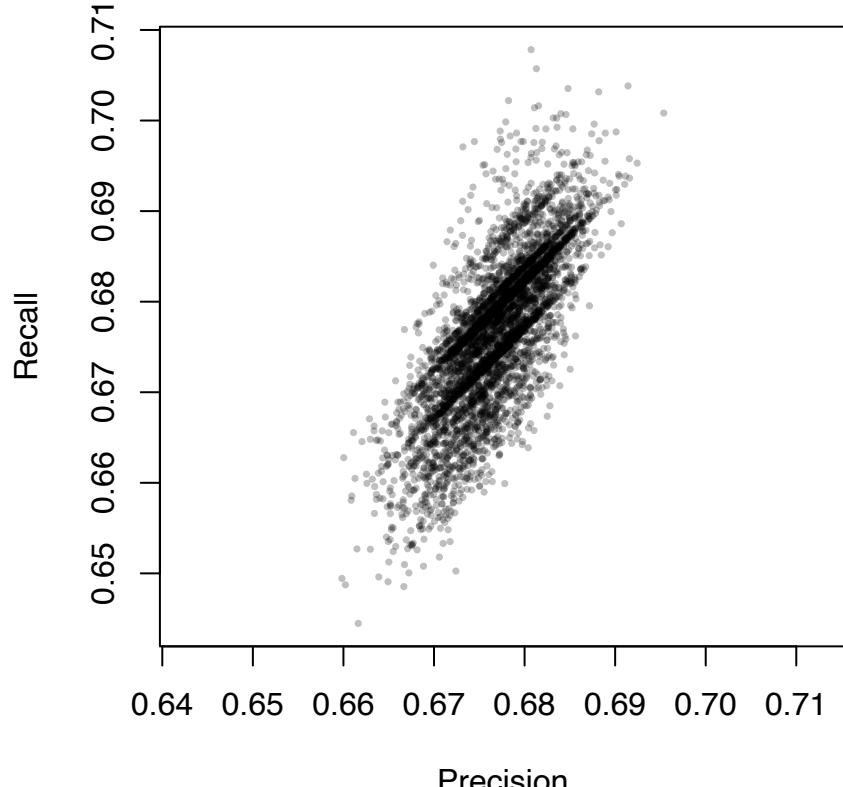
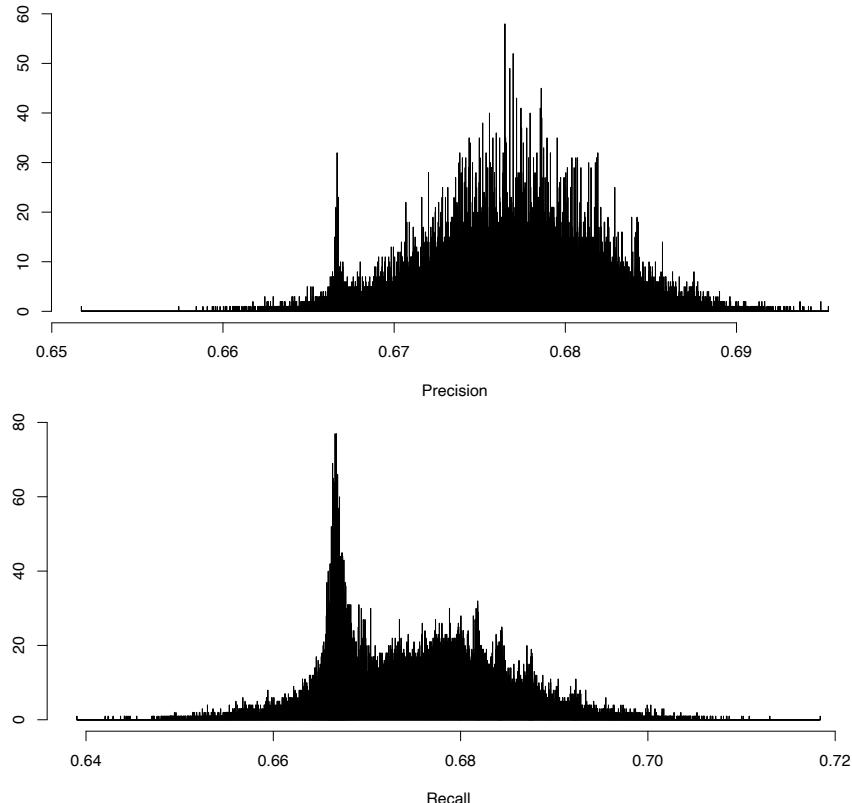




Distributions of precision/recall? $n = 10^4$

($n_{\text{sim}} = 50,000$)

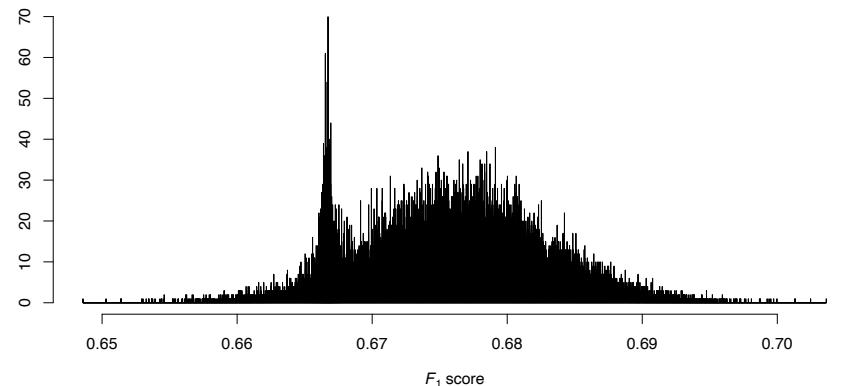
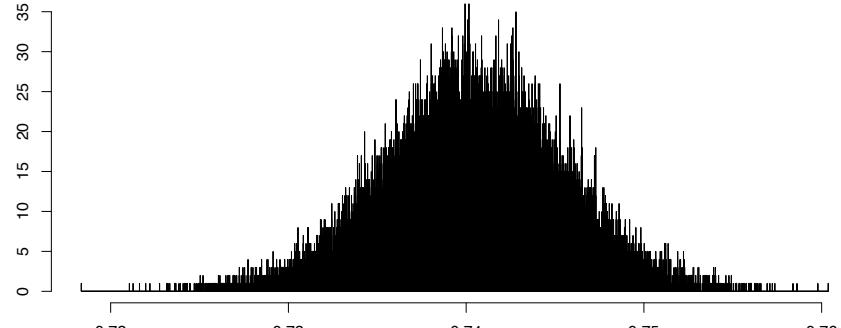
- Introduction
- Sampling frame and measurement
- "Prediction" vs. causality
- Model metrics as estimators**
- Cultural issues
- Lessons from other fields
- Summary and conclusion
- References
- Appendix: Simulation code





Distributions of AUC/ F_1 ? $n = 10^4$ ($n_{\text{sim}} = 50,000$)

- 95% tolerance interval is tighter, but still larger than state-of-the-art improvements
- Also... are precision, recall, and F_1 mixtures?? That is really weird!
 - **How has nobody looked at this??**
The distribution of estimators is stats theory 101!
 - (The problem persists for various seeds. But maybe I made a mistake?)
- Conclusion: even for large sample size, a simple DGP, and a "true" model, the distribution of common metrics is not simple





Take-aways: We should study asymptotic distributions of metrics, and use them!

- Can somebody please find the distributions of ML model success metrics? (I started to try, via joint distribution of TP, FP, FN, TN as a multidimensional [3+1 dimensions] binomial, and then taking ratios of marginals, but it's a lot of algebra)
- With distributions, we could find asymptotic confidence intervals, and conduct significance testing of model results
 - Yes, p -values and hypothesis testing have done enormous damage, **but ignoring variance might be worse**
 - Also, we should probably start doing **power calculations** in ML
- Maybe, when studying asymptotic distributions, we'll find sufficient statistics for model success (like the parameters of a multivariate binomial) and good estimators thereof

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

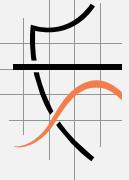
Cultural issues

Lessons from other fields

Summary and conclusion

References

Appendix:
Simulation code



Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

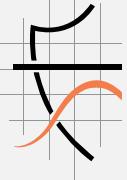
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

Cultural issues



Narrow technical training

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- Phil Agre (1997):
 - "My college did not require me to take many humanities courses, or learn to write in a professional register, and so I arrived in graduate school at MIT with little genuine knowledge beyond math and computers. This realization hit me with great force halfway through my first year of graduate school..."
 - "I was unable to turn to other, nontechnical fields for inspiration... The problem was not exactly that I could not understand the vocabulary, but that I insisted on trying to read everything as a narration of the workings of a mechanism."
- Study design and measurement still partially fall under "technical" knowledge; the problem is far more profound



Paradigms of inquiry

	Issue	Positivism	Post-positivism	Critical theory et al.	Constructivism	Participatory
Introduction						
Sampling frame and measurement	Ontology	Naïve realism: Reality independent of and prior to human conception of it, apprehensible	Critical realism: Reality independent of and prior to human conception, but imperfectly and approx. apprehensible	Disenchantment theory: reality is secret/hidden, shaped by power structures and solidified over time	Relativism: multiple realities, constructed in history through social processes	Participative: multiple realities, co-constructed through interactions between specific people and environments
"Prediction" vs. causality	Epistemology	Reality knowable. Findings are singular, neutral, perspective-independent, atemporal, universally true	Findings provisionally true; multiple descriptions can be valid but are probably equivalent; findings can be affected/distorted by social + cultural factors	How we come to know something, or who knows it, matters for how meaningful it is	Relativistic: no neutral perspective to adjudicate competing claims	We come to know things, create new understandings, & transform world by involving other people in process of inquiry
Model metrics as estimators						
Cultural issues	Methodology	Hypotheses can be verified as true. Quant methods, math.	Falsification of hypotheses; primacy of quant, but some qual and mixed methods	Dialogic (conversation + debate) or dialectical ($\text{thesis}_1 \rightarrow \text{antithesis}_1 \rightarrow \text{synthesis}_2 := \text{thesis}_2\dots$)	Dialectical, or exegetical (reading between the lines")	Collaborative, action-focused; flattening hierarchies; engaging in self- and collective reflection, action
Lessons from other fields						
Summary and conclusion	Axiology	Quant knowledge-holders have access to truth, and responsibility from it	Quant knowledge valuable but can be distorted; qual can help find and correct	Marginalization provides unique insights, knowledge of marginalized valuable	Understanding construction is valuable; value relative to given perspective	Reflexivity, co-created knowledge, and non-western ways of knowing are valuable and combat erasure and dehumanization
References						
Appendix: Simulation code						

Malik & Malik (2021), via Guba and Lincoln (2005)



Ways of understanding a person

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

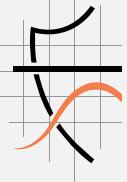
Summary and
conclusion

References

Appendix:
Simulation
code

	As a case (quant)	In narrative (qual)
Context/ circumstance	Stripped away	Key
Mental states	Absent (for the most part)	Crucial; constitutive
Relevant features	Determined in advance	Emergent
Orientation to time	Atemporal	Chronological
Ordering of features	Unimportant	Meaningful
Other actors	Invisible	Often present
Causal logic	Mathematical	Theoretical
Boost predictive validity	Add cases	Know person better

Slide from Barbara Kiviat (work in progress), based on "Bowker and Star 2000; Bruner 1986; Desrosières 1998; Espeland 1998; Espeland and Stevens 1998, 2008; Fourcade and Healy 2017; Hacking 1990; Porter 1994, 1995; Ricouer 1998; White 1980, 1984". I would add: Patton 2005; Abbott 1988



Why this matters: it's why we expect generalization

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

Cultural issues

Lessons from other fields

Summary and conclusion

References

Appendix:
Simulation code

- We expect that models are picking up on signal, not noise
 - Statistics makes the assumption that we can treat the world as made up of entities that are distinct but are realizations of an underlying process. Machine learning shares this assumption, even if it is not explicit about it (e.g., theory about convergence to the "oracle predictor" rather than about convergence to a "true" parameter)
- If we define the "signal" as what is invariant, then failures of generalizability means we've failed to find the underlying regularity
- But is there really aggregate regularity? Or only *narrative*, if any?
 - E.g., Twitter and elections (Gayo-Avello, 2012)
 - Note: one explanation for stats working is that it *imposes* regularity



“What are we even doing?”

Introduction

Sampling
frame and
measurement

“Prediction”
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- “If science isn’t ‘true’, then what are we even doing? We might as well be doing English literature, or art criticism!”
 - Intellectual supremacy is probably a bad reason for doing science
- At least for the social world, I am skeptical of attempts to find underlying regularity in the [social] world as cases; both because only trivial things can have universal aggregate regularity, and because attempts to find social regularity can end up imposing it
 - But neither can I imagine our civilization without the use of summary statistics for management, planning, and allocation...



Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

**Lessons from
other fields**

Summary and
conclusion

References

Appendix:
Simulation
code

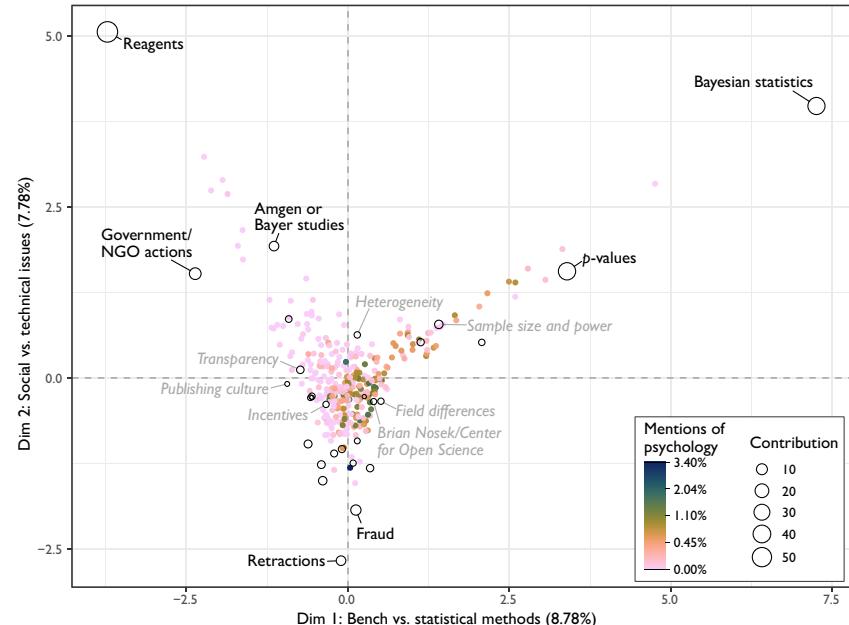
Lessons from other fields



Proximate causes? Concerns?

Introduction
Sampling frame and measurement
"Prediction" vs. causality
Model metrics as estimators
Cultural issues
Lessons from other fields
Summary and conclusion
References
Appendix: Simulation code

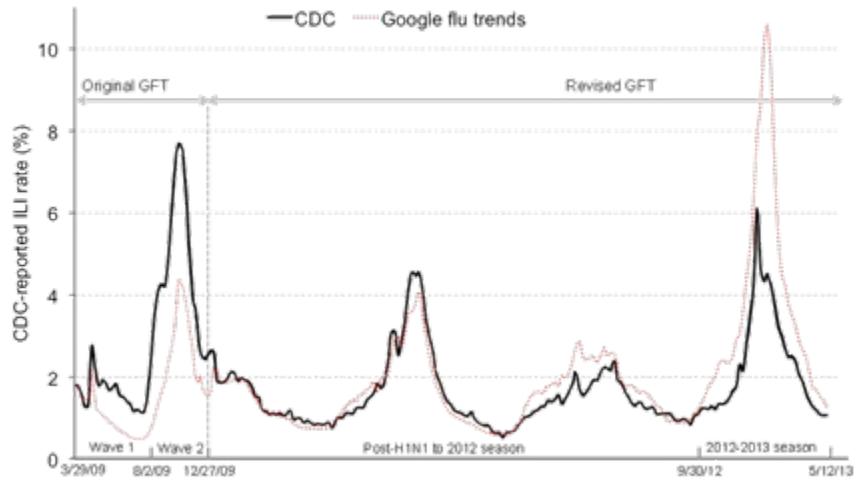
- Attention in 2011/2012 in both psychology and biomedicine, rapidly led to major policy initiatives
- What caused? In psych: one researcher admitting to a decade of (deliberate) fraud that went undiscovered + the "ESP" study + decades of concern → "soul searching"
- Psychologists and collaborators thought the crisis/problems were worse there: stat ignorance ("methodolatry")? Objects of study more variable? Bad incentives?
- But discourse in biomedicine was very similar: clustering doesn't separate fields (Nelson, Chung, Ichikawa, & Malik, 2022)





Take-aways

- Machine learning is not necessarily special in having a crisis, or worse than other fields
- In other fields, **dramatic failures (rather than long-standing concerns) precipitated an experience of "crisis"**
- Machine learning had a dramatic failure in Google Flu Trends, which Lazer (2014) called "big data's 'Dewey Defeats Truman' moment"
 - But that didn't prompt much reform
 - Epic Sepsis model? Also a flash in pan





What can we expect? What should we want?

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

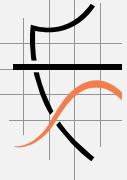
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- We probably won't get large-scale reform attempts until we enter a crisis. Crisis requires dramatic failure and attention
 - If we really want reform, maybe we should *want* a crisis, and try to precipitate one...
- Hype about claimed success are probably enough to prevent getting around high-profile failures for some time
 - Consequences of a crisis? Maybe loss of legitimacy and funding (another "AI winter"): but also, if hype is sufficient there is a niche for "reformers" (Nelson, 2018) who preserve legitimacy and funding
- The fundamental sources of the problem: yes, methods and incentives, both of which we can and should improve
- A different solution: I don't think replication should be the measure of science, such that failures of replication shouldn't be that big a deal



Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

**Summary and
conclusion**

References

Appendix:
Simulation
code

Summary and conclusion



Contextual issues

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

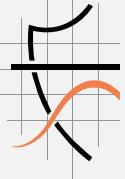
Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

- Maybe irreproducibility is just an artifact of a positivist commitment: if we change our understanding of what science is, should be, and could be to something far humbler, then reproducibility wouldn't be a problem
- Based on prior fields, and on current hype, we shouldn't expect reform without a crisis and we shouldn't expect a crisis anytime soon
- But there are still methodological steps we can and should take. These will hopefully fix any problems for ML applications to physical sciences and some engineering



Methodological issues

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code



ML models will only generalize insofar as the **data are representative**



Selection on the dependent variable is not something we can do in application



If the **underlying measurements** are not consistent, the model can also fail to generalize



Point estimates of **model metrics** don't give possible **variability** even with the same population



Dependencies cause a form of **leakage**



Unless models give unbiased estimates of partial correlation, **causal shifts** will make them invalid



Suggested fixes

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code



Gather representative data and/or make more limited claims



Include weak signal observations, rather than filter them out



Use a **measurement model** (for the response), or at least consider validity and reliability



Get **confidence intervals** around all measures of model success, and **study asymptotics**



Split data by dependencies (temporal block CV, leave-one-subject-out CV, network CV, etc.)



Change language to temper expectations, and **sometimes, pursue causality**



References (1/3)

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

Cultural issues

Lessons from other fields

Summary and conclusion

References

Appendix:
Simulation code

Arceneaux, Kevin, Alan S. Gerber, and Donald P. Green. 2010. 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 (2):256-282. <https://dx.doi.org/10.1177/0049124110378098>

Borgatti, Steve. 2012. Types of validity. BA 762: Research Methods. Gatton College of Business & Engineering, University of Kentucky.
<https://sites.google.com/site/ba762researchmethods/reference/handouts/types-of-validity>

Box, George E. P. 1979. *Robustness in the strategy of scientific model building*. Technical Report #1954. Mathematics Research Center, University of Wisconsin-Madison.

Breiman, Leo. 2001. Statistical modeling: The two cultures. *Statistical Science* 16 (3): 199–231.
<https://dx.doi.org/10.1214/ss/1009213726>

Bryman, Alan. 1988. *Quantity and quality in social research*. Routledge. <https://doi.org/10.4324/9780203410028>

Cardoso, Fatima, Laura J. van't Veer, Jan Bogaerts, Leen Slaets, Giuseppe Viale, Suzette Delaloge, Jean-Yves Pierga, Etienne Brain, Sylvain Causeret, Mauro DeLorenzi, Annuska M. Glas,

Vassilis Golfinopoulos, Theodora Goulioti, Susan Knox, Erika Matos, Bart Meulemans, Peter A. Neijenhuis, Ulrike Nitz, Rodolfo Passalacqua, Peter Ravdin, Isabel T. Rubio, Mahasti Saghatelian, Tineke J. Smilde, Christos Sotiriou, Lisette Stork, Carolyn Straehle, Geraldine Thomas, Alastair M. Thompson, Jacobus M. van der Hoeven, Peter Vuylsteke, René Bernards, Konstantinos Tryfonidis, Emiel Rutgers, and Martine Piccart. 2016. 70-gene signature as an aid to treatment decisions in early-stage breast cancer. *New England Journal of Medicine* 375 (8): 717-729. <https://dx.doi.org/10.1056/NEJMoa1602253>

Dwork, Cynthia, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth. 2015. The reusable holdout: Preserving validity in adaptive data analysis." *Science* 349 (6248): 636-638. <https://dx.doi.org/10.1126/science.aaa9375>

Efron, Bradley. 2004. The estimation of prediction error: Covariance penalties and cross-validation. *Journal of the American Statistical Association* 99 (467): 619–632.
<https://dx.doi.org/10.1198/016214504000000692>

Gilbert, Nigel, and Klaus Troitzsch. 2005. *Simulation for the social scientist*. 2nd edition. Open University Press.



References (2/3)

Introduction

Sampling frame and measurement

"Prediction" vs. causality

Model metrics as estimators

Cultural issues

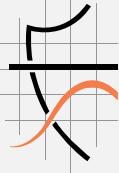
Lessons from other fields

Summary and conclusion

References

Appendix:
Simulation code

- Ginsberg, Jeremy, Matthew H. Mohebbi, Rajan S. Patel, Lynnette Brammer, Mark S. Smolinski, and Larry Brilliant. 2009. Detecting influenza epidemics using search engine query data. *Nature* 457: 1012-1015. <https://dx.doi.org/10.1038/nature07634>
- Hoadley, Bruce. 2001. [Statistical modeling: The two cultures]: Comment. *Statistical Science* 16 (3): 220-224.
- Imbens, Guido W., and Thomas Lemieux. 2008. Regression discontinuity designs: A guide to practice. *Journal of Econometrics*, 142 (2): 615-635.
<https://doi.org/10.1016/j.jeconom.2007.05.001>
- Jacobs, Abigail Z., and Hanna Wallach. 2021. Measurement and fairness. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (FAccT '21), 375-85.
<https://doi.org/10.1145/3442188.3445901>
- Jones, Matthew L. 2015. How we became instrumentalists (again): Data positivism since World War II. *Historical Studies in the Natural Sciences* 48 (5): 673-684.
<https://dx.doi.org/10.1525/hsns.2018.48.5.673>
- Kass, Robert E. 2011. Statistical inference: The big picture. *Statistical Science* 26 (1): 1-9. <https://dx.doi.org/10.1214/10-STS337>
- Keeling, Matt J., and Pejman Rohani. 2008. *Modeling infectious diseases in humans and animals*. Princeton University Press.
<https://doi.org/10.1515/9781400841035>
- Kleinberg, Jon, Jens Ludwig, Sendhil Mullainathan, and Ziad Obermeyer. 2015. Prediction policy problems. *American Economic Review* 105 (5): 491-495.
<https://dx.doi.org/10.1257/aer.p20151023>
- Koren, Yehuda. 2009. Collaborative filtering with temporal dynamics.
- Lazer, David, Ryan Kennedy, Gary King, and Alessandro Vespignani. The parable of Google Flu: Traps in big data analysis. *Science* 343 (6176): 1203-1205.
<https://dx.doi.org/10.1126/science.1248506>
- Malik, Maya and Momin M. Malik. 2021. Critical technical awakenings. *Journal of Social Computing* 2 (4): 365-384.
<https://doi.org/10.23919/JSC.2021.0035>
- Mullainathan, Sendhil and Jann Spiess. 2017. Machine learning: An applied econometric approach. *Journal of Economic Perspectives* 31 (2): 87-106.
<https://dx.doi.org/10.1257/jep.31.2.87>
- Olteanu, Alexandra, Carlos Castillo, Fernando Diaz, and Emre Kiciman. 2019. Social data: Biases, methodological pitfalls, and ethical boundaries. *Frontiers in Big Data* 2(13): 1-33.
<https://doi.org/10.3389/fdata.2019.00013>



References (3/3)

Introduction
Sampling frame and measurement
"Prediction" vs. causality
Model metrics as estimators
Cultural issues
Lessons from other fields
Summary and conclusion
References
Appendix: Simulation code

- Olteanu, Alexandra, Carlos Castillo, Fernando Diaz, and Emre Kiciman. 2019. Social data: Biases, methodological pitfalls, and ethical boundaries. *Frontiers in Big Data* 2(13): 1–33. <https://doi.org/10.3389/fdata.2019.00013>
- Park, Greg. 2012. The dangers of overfitting: A Kaggle postmortem. <http://qrgpark.io/blog/Kaggle-Psychopathy-Postmortem/>
- Rescher, Nicholas. 1998. *Predicting the future: An introduction to the Teele theory of forecasting*. State University of New York Press.
- Richardson, Eugene T., Momin M. Malik, William A. Darity, Jr., A. Kirsten Mullen, Michelle E. Morse, Maya Malik, Adia Benton, Mary T. Bassett, Paul E. Farmer, Lee Worden, and James Holland Jones. 2021. Reparations for Black American descendants of persons enslaved in the U.S. and their potential impact on SARS-CoV-2 transmission. *Social Science & Medicine* 276: 113741. <https://doi.org/10.1016/j.socscimed.2021.113741>
- Santillana, Mauricio, Wendong Zhang, Benjamin M. Althouse, and John W. Ayers. 2014. What can digital disease detection learn from (an external revision to) Google Flu Trends? *American Journal of Preventive Medicine* 47 (3): 341–347. <http://dx.doi.org/10.1016/j.amepre.2014.05.020>
- Savage, Mike, and Roger Burrows. 2007. The coming crisis of empirical sociology. *Sociology* 41 (5): 885–899. <https://doi.org/10.1177/0038038507080443>
- Shmueli, Galit. 2010. To explain or to predict? *Statistical Science* 25 (3): 289–310. <https://dx.doi.org/10.1214/10-STS330>
- Tasse, Dan, 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), 250–259.
- Teele, Dawn Langan. 2014. *Field experiments and their critics: Essays on the uses and abuses of experimentation in the social sciences*. Yale University Press.
- van't Veer, Laura J., Hongyue Dai, Marc J. van de Vijver, Yudong D. He, Augustinus A. M. Hart, Mao Mao, Hans L. Peterse, Karin van der Kooy, Matthew J. Marton, Anke T. Witteveen, George J. Schreiber, Ron M. Kerkhoven, Chris Roberts, Peter S. Linsley, René Bernards, and Stephen H. Friend. Gene expression profiling predicts clinical outcome of breast cancer. *Nature* 415, no. 6871 (2002): 530–536. <https://dx.doi.org/10.1038/415530a>
- Zagheni, Emilio, and Ingmar Weber. 2015. Demographic research with non-representative internet data. *International Journal of Manpower* 36 (1): 13–25. <https://doi.org/10.1108/IJM-12-2014-0261>



Appendix: Simulation code

Introduction

Sampling
frame and
measurement

"Prediction"
vs. causality

Model metrics
as estimators

Cultural issues

Lessons from
other fields

Summary and
conclusion

References

Appendix:
Simulation
code

```
library(ModelMetrics)

# Rename for convenience
logistic <- function(x) plogis(x)
logit <- function(p) qlogis(p)

set.seed(20220728)
nsim <- 50000
results <- data.frame(accuracy = rep(NA, nsim),
                      ppv = rep(NA, nsim),
                      tp = rep(NA, nsim),
                      tn = rep(NA, nsim),
                      fp = rep(NA, nsim),
                      fn = rep(NA, nsim),
                      tpr = rep(NA, nsim),
                      tnr = rep(NA, nsim),
                      auc = rep(NA, nsim),
                      f1score = rep(NA, nsim))

# Either run with 97 or 101 (small sample size):
# these are prime number close to 100, so
# that accuracy and other fractions divided by
# a prime denominator), or 10k (large sample
# size)

# n <- 97
# n <- 101
n <- 10000

# Draw X once ("fixed X" setting), then draw a new Y
# each simulation run, y ~ bernoulli(logistic(x))
x <- rnorm(n = n, mean = 0, sd = 1)

for (i in 1:nsim) {
  y <- rbinom(n = n, size = 1, prob = logistic(x))
  glml <- glm(y ~ x, family = "binomial")
  results$accuracy[i] <- mean(y==(predict.glm(glml, type = "response")>.5))
  results$ppv[i] <- ppv(y, predict.glm(glml, type = "response")) # Precision
  results$tp[i] <- sum(y==1 & (predict.glm(glml, type = "response")>=.5))
  results$tn[i] <- sum(y==0 & (predict.glm(glml, type = "response")<.5))
  results$fp[i] <- sum(y==0 & (predict.glm(glml, type = "response")>=.5))
  results$fn[i] <- sum(y==1 & (predict.glm(glml, type = "response")<.5))
  results$tpr[i] <- tpr(y, predict.glm(glml, type = "response")) # Recall
  results$tnr[i] <- tnr(y, predict.glm(glml, type = "response")) # Specificity
  results$auc[i] <- auc(y, predict.glm(glml, type = "response"))
  results$f1score[i] <- f1Score(y, predict.glm(glml, type = "response"))
  if (i%%1000==0) {print(i)}
}
```