



# Critical Approaches to Machine Learning

**Momin M. Malik**

Sunday, 27 March 2022

ICQCM 2022 Summit

Baltimore, MD

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Overview

- Background/review
- Scope; review of some critiques of machine learning
- Examples of work that take, or enable, a critical approach
- Possibilities



# Goals

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

- To separate out possibilities “external” to ML to those “internal” to ML (while also critiquing, and blurring, that boundary)
- To give neat examples of ML
- To connect some specific technical constructs in ML to possibilities for critique
- To get a conversation/sharing going!



Overview

**Background/  
review**

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Background/review



# Building on two remote sessions

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

- “Defining Critical Quantitative and Computational Methodologies.” William T. Grant AQC SCHOLARS Virtual Seminar Series, 27 May 2021.  
<https://www.mominmalik.com/icqcm2021.pdf>
- “Computational Approaches III: Applications.” ICQCM 2021 Seminar Series, 22 July 2021.  
<https://www.mominmalik.com/icqcm2021b.pdf>
- If you didn’t know and/or didn’t attend those sessions: this builds on them, but doesn’t require them. Also feel free to check out the slides!

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# “Defining QCM” key points: What

- Incorporating ready-made quant methods into a critical approach is most straightforward
  - “Minimal” critical QCM: quant demonstrations of disparities *that links to theory* about the source of those disparities (e.g., white supremacy, dehumanization)
- But more intellectually interesting for me is integrating the logic of modeling with the logic of critical theory at a fundamental level
- Much harder—requiring dual training—but a rich intellectual project
  - On the other hand, maybe useless practically, and the “minimal” version of critical QCM is most useful and important



# “Defining QCM” key points: Why

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

- Strategic quantification/strategic modeling (after Spivak’s “strategic essentialism”) to demonstrate inequality?  
(Strategic positivism; Wyly, 2009)
  - Rhetorical use: convince power-brokers?
- “Counterhegemonic modeling” (Richardson 2020): modeling ironically to reveal the absurdity of modeling?
- Alternatively: just because quantification is currently associated with power does not mean it is essentially so. Qualitative inquiry can be just as or more oppressive, it just isn’t currently in power

# "Computational Approaches" key points

Overview

Background/  
review

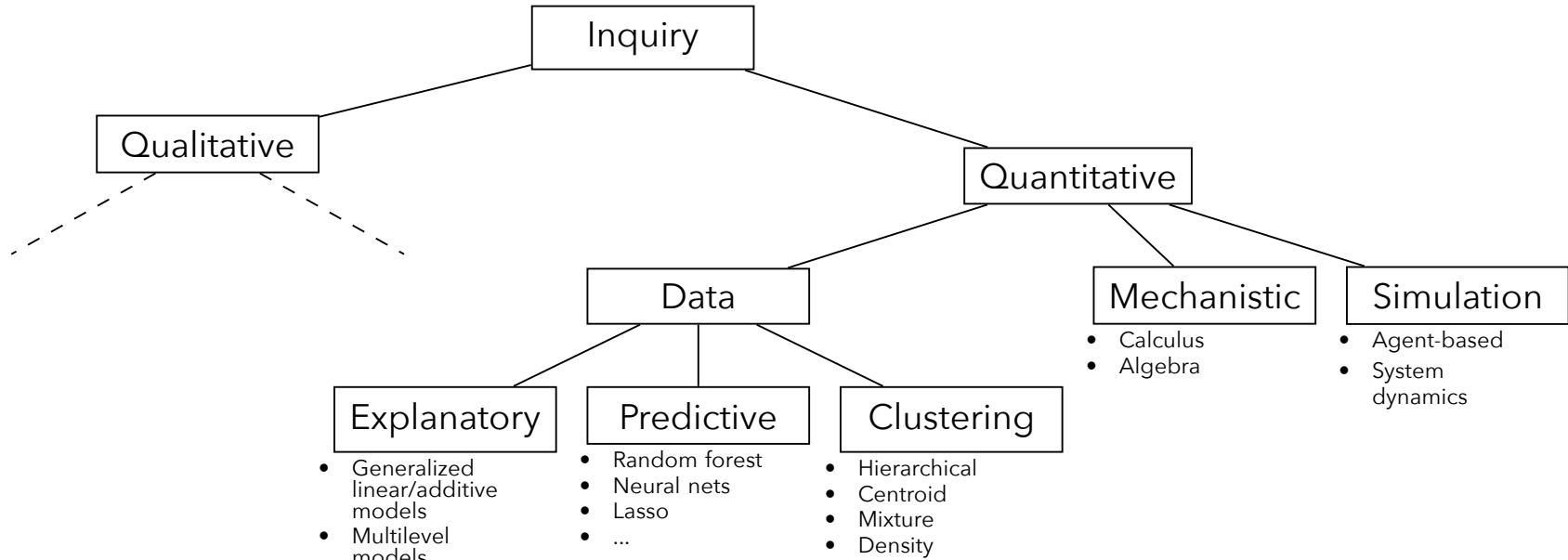
Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References



Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

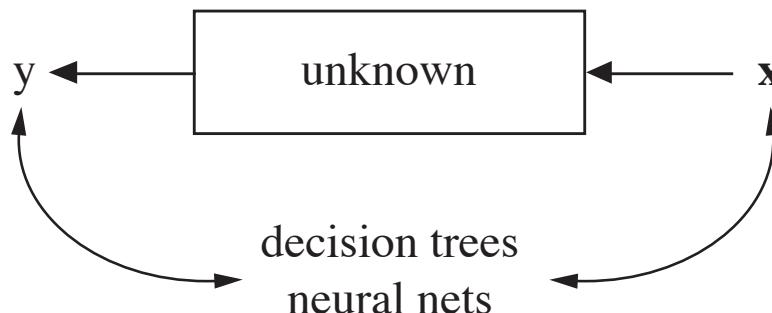
References

# "Computational Approaches" key points

## Statistics:



## Machine learning:



Machine learning: An instrumental use of statistical correlations to *mimic* the output of a target process, rather than understand the *relationship* between inputs and outputs. Involves finding expressions that maximize correlation.

Breiman 2001. See also Jones 2018.



# "Computational Approaches" key points

**"Source subject": Marquese Scott**

## Everybody Dance Now

Motion Retargeting Video Subjects

Caroline Chan, Shiry Ginosar, Tinghui Zhou, Alexei A. Efros

UC Berkeley

Caroline Chan, "Everybody Dance Now: Motion Retargeting Video Subjects."  
<https://youtu.be/PCBTZh41Ris>



Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Scope; Review of some critiques



# ML is bad for making claims about the world!!

- *The best-fitting (most accurate\*) model does not necessarily reflect how the world works*
- This has been shocking in statistics for decades (Stein's paradox, Leo Breiman's "two cultures"), but little known outside
- Why: one reason is the "bias-variance tradeoff"
  - Even when available, the "true" covariates may be noisy, in which case proxies (or even just going with the mean) sometimes does better
- Another reason: narrowing in to get one causal relationship "correct" might require sacrificing the rest of the model
- So: we can use correlations to "predict" without "explaining" (knowing causality)!

\* Or other relevant metric of success

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

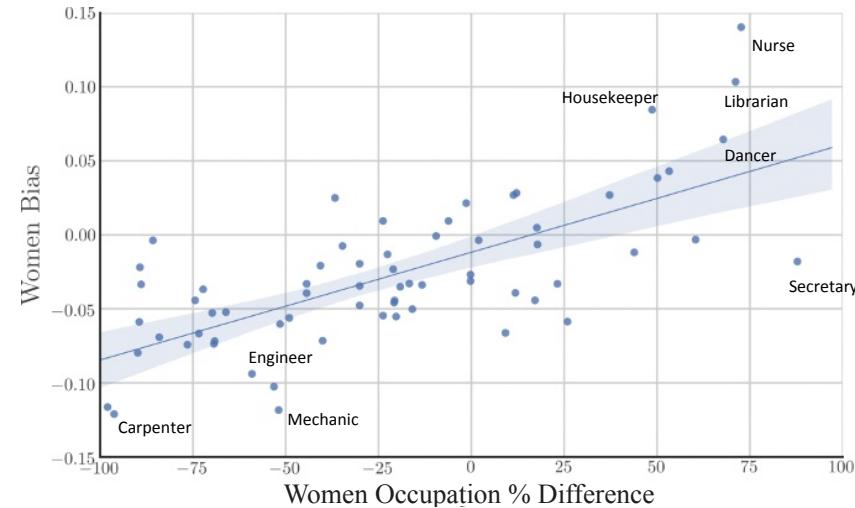
Possibilities

Conclusion  
and  
discussion

References

# Machine learning to “prove”?

- Unlike statistics, machine learning isn't well-suited for making claims about the way the world works
- Maybe incidental outputs of models are revealing (e.g., “Word embeddings quantify 100 years of gender and ethnic stereotypes”)... but that is indirect, and reifies ML





Overview

Background/  
review

**Scope;**  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Critiques of ML?

- Safiya Noble
- Ruha Benjamin
- Meredith Broussard
- Virginia Eubanks
- Cathy O'Neil
- Adrian Mackenzie
- Dan McQuillan
- Matteo Pasquinelli
- Yarden Katz
- *Data Feminism*
- Coalition for Critical Technology

Overview

Background/  
review

**Scope;**  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Applied STS?

- Incorporating science studies lessons (e.g., “responsible research and innovation”) into data science education
  - “Integrating FATE/Critical Data Studies into Data Science Curricula: Where are we going and how do we get there?”
  - “Critique and contribute: A practice-based framework for improving critical data studies and data science”
  - “‘You Social Scientists Love Mind Games’: Experimenting in the ‘divide’ between data science and critical algorithm studies”

# Can critique and quantification mix?

Overview

Background/  
review

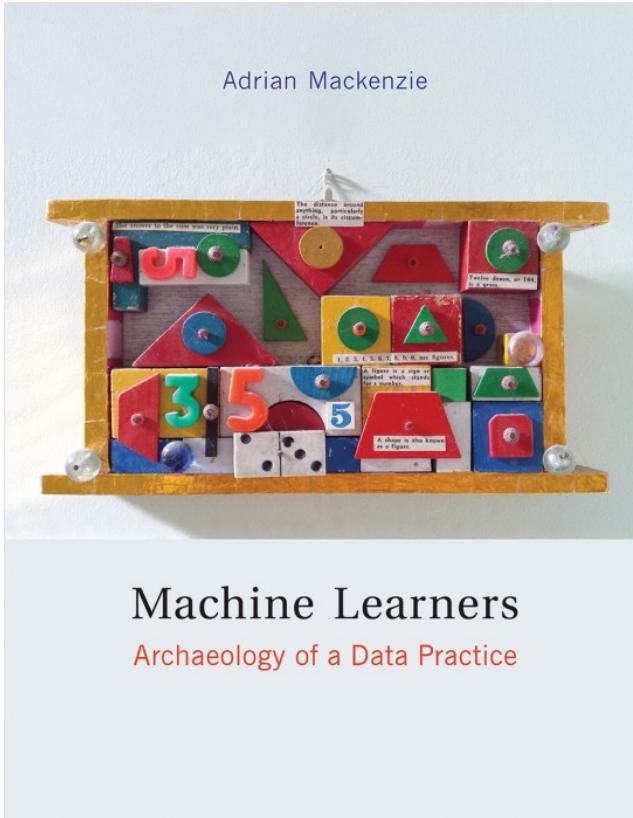
**Scope;**  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References



We might also approach the epistemic fault line in machine learning topologically. More than a decade ago, the cultural theorist Brian Massumi wrote that "the space of experience is really, literally, physically a topological hyperspace of transformation" (Massumi 2002, 184). Much earlier, Gilles Deleuze had conceptualized Michel Foucault's philosophy as a topology, or "thought of the outside" (Deleuze 1988b), as a set of movements that sought to map the diagrams that generated a "kind of reality, a new model of truth" (Deleuze 1988b, 35). More recently, this topological thinking has been extended and developed by Celia Lury among others. In "The Becoming Topological of Culture," Lury, Luciana Parisi, and Tiziana Terranova suggest that "a new rationality is emerging: the moving ratio of a topological culture" (Lury, Parisi, and Terranova 2012, 4). In this new rationality, practices of ordering, modeling, networking, and mapping co-constitute culture, technology, and science (Lury, Parisi, and Terranova 2012, 5). At the core of this new rationality, however, lies a new ordering of continuity. The "ordering of continuity," Lury, Parisi, and Terranova propose, takes shape "in practices of sorting, naming, numbering, comparing, listing, and calculating" (4). The phrase "ordering of continuity" is interesting because we don't normally

## Vectorization and Its Consequences

65

and Andrew Ng advocate returning often to equations). The mainstay of statistics, the linear regression model, usually appears in a more or less algebraic form:

$$\hat{Y} = \hat{\beta}_0 + \sum_{j=1}^p X_j \hat{\beta}_j \quad (3.1)$$

$$\hat{Y} = X_T \hat{\beta} \quad (3.2)$$

Equations 3.1 and 3.2 express a plane (or hyperplane) in increasingly diagrammatic abstraction. The possibility of diagramming a high-dimensional space derives largely from linear algebra. Reading equation 3.1 from left to right, the expression  $\hat{Y}$  already



# Artistic approaches?

## Overview

## Background/ review

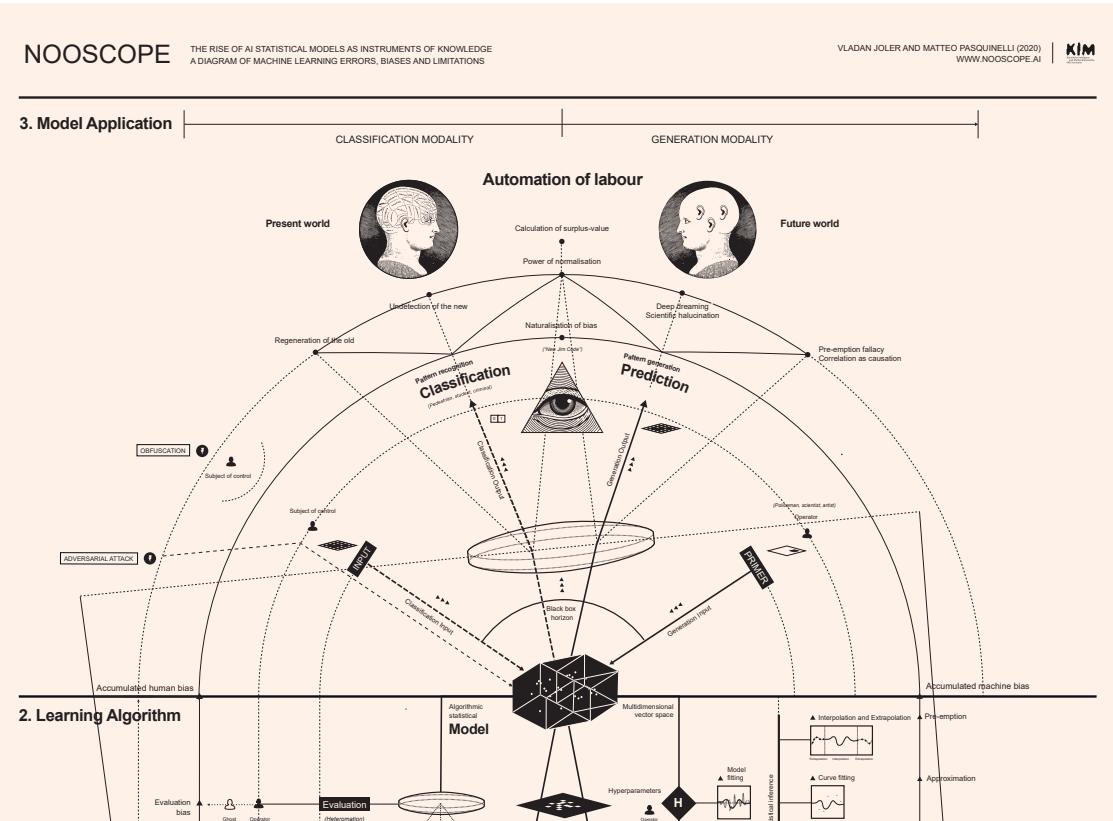
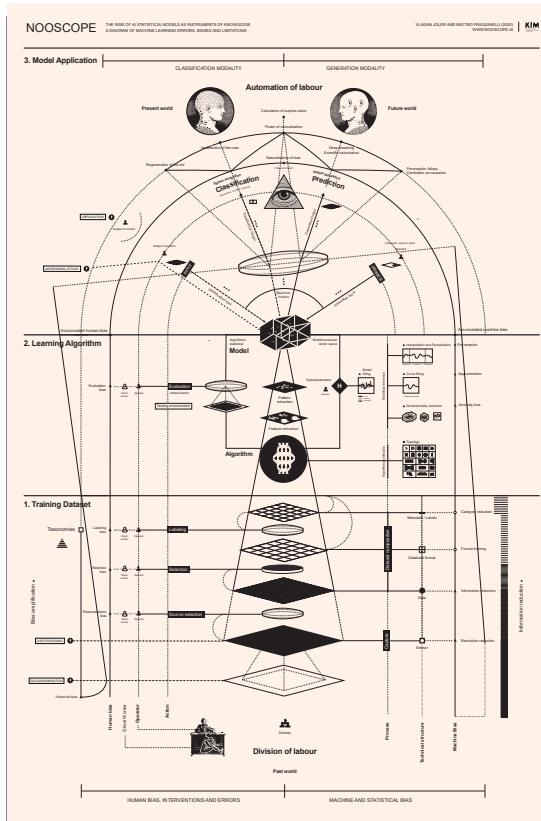
## Scope; review of some critiques

Examples of work that take, or enable critique

## Possibilities

## Conclusion and discussion

## References





# General methodological critique?

Overview

Background/  
review

**Scope;**  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

- “Troubling Trends in Machine Learning Scholarship”
- “Machine learning that matters”
- “Reliance on metrics is a fundamental challenge for AI”
- “Underspecification presents challenges for credibility in modern machine learning”



# A Survey of Papers I Think Are Neat/Important And Why

**Momin M. Malik**

Sunday, 27 March 2022

ICQCM 2022 Summit

Baltimore, MD

[Overview](#)
[Background/  
review](#)
[Scope;  
review of  
some  
critiques](#)
[Examples of  
work that  
take, or  
enable  
critique](#)
[Possibilities](#)
[Conclusion  
and  
discussion](#)
[References](#)

# Auditing

- “Gender shades:  
Intersectional accuracy  
Disparities in commercial  
gender classification”
- Follow-up work: “Actionable  
auditing: Investigating the  
impact of publicly naming  
biased performance results  
of commercial AI products”,  
“Lessons from archives:  
strategies for collecting  
sociocultural data in  
machine learning”

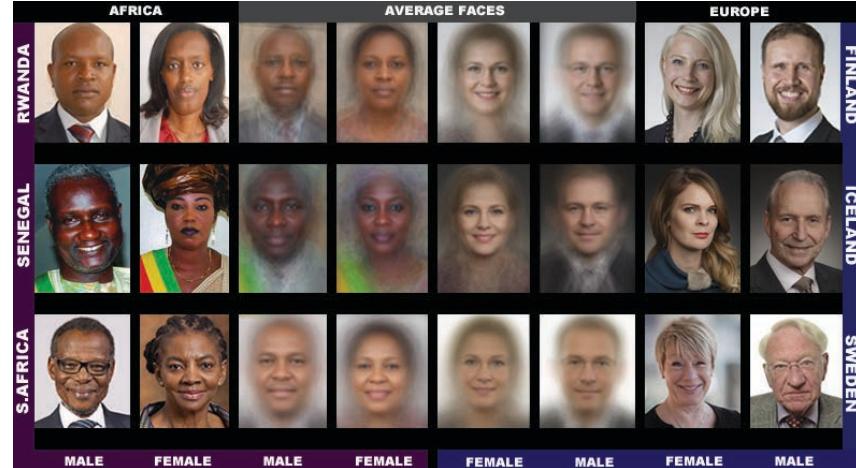


Figure 1: Example images and average faces from the new Pilot Parliaments Benchmark (PPB). As the examples show, the images are constrained with relatively little variation in pose. The subjects are composed of male and female parliamentarians from 6 countries. On average, Senegalese subjects are the darkest skinned while those from Finland and Iceland are the lightest skinned.

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# ML as scalable measurement?

- “Constructing a visual dataset to study the effects of spatial apartheid in South Africa”
- (See also: “Measuring urban social diversity using interconnected geo-social networks”

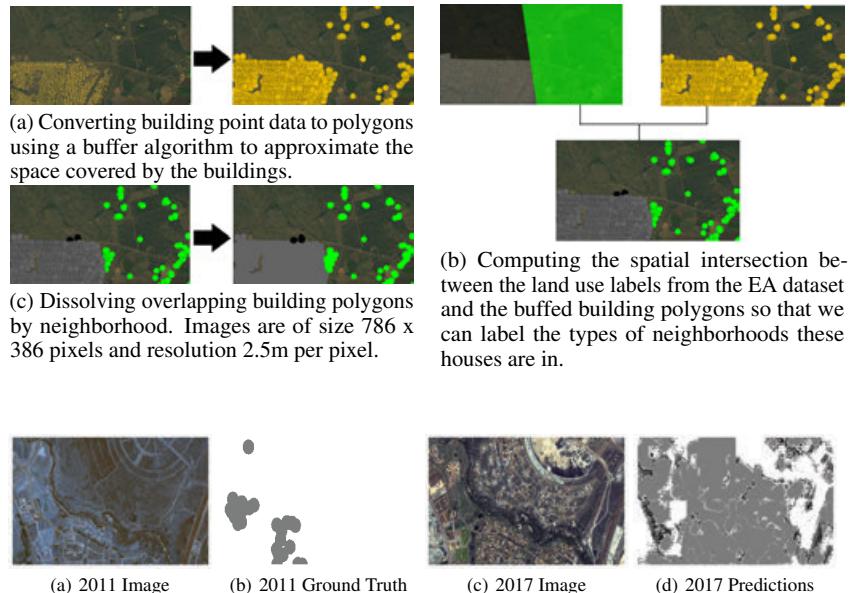
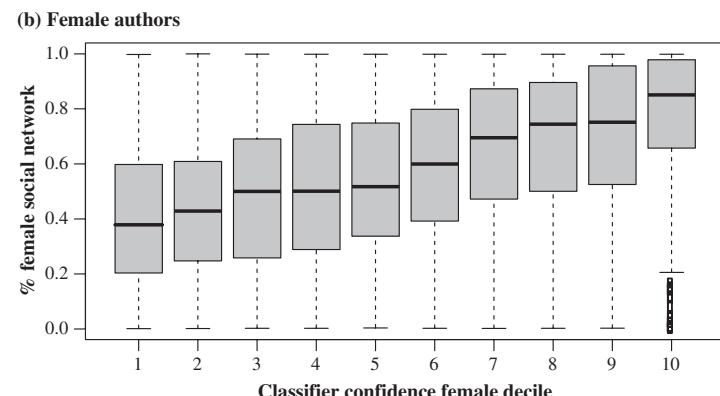
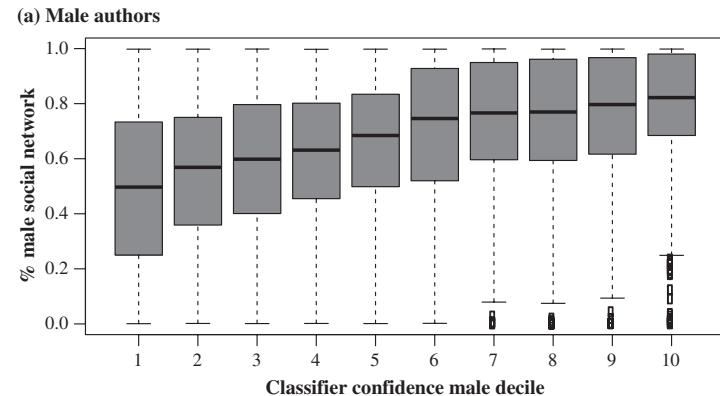


Figure 7: Examples of the change detected between 2011 and 2017 images in a wealthy neighborhood near a big mall. Dark gray: Wealthy Neighborhood, White: Background.

[Overview](#)
[Background/  
review](#)
[Scope;  
review of  
some  
critiques](#)
[Examples of  
work that  
take, or  
enable  
critique](#)
[Possibilities](#)
[Conclusion  
and  
discussion](#)
[References](#)

# Qualitative follow-up?

- “Gender identity and lexical variation in social media”
- Classified along a gender binary in line with existing work, but then looks at people who are “misclassified” qualitatively, in recognition of gender as a performance
- (See also: “Gender recognition or gender reductionism? The social implications of automatic gender recognition systems”, “The Misgendering Machines: Trans/HCI Implications of Automatic Gender Recognition”)



# Qualitative evaluation?

- “Predictive learning analytics ‘at scale’: Towards guidelines to successful implementation in higher education based on the case of the Open University UK”
- (See also: “A large-scale implementation of predictive learning analytics in higher education: The teachers’ role and perspective”

**Table 1:** Themes from Interviews with 20 Education Managers (Organized by Faculty)

Themes	RQs	Science faculty	Business Faculty	Teaching and Learning Unit	Student engagement and support	Tuition delivery	Data professionals
General perceptions	RQ1	Positive	Positive	Positive	Positive	Positive	Positive
Perceived challenges	RQ2	Lack of evaluation (Rec. 1)*; Lack of understanding as to how to use PLA (Rec. 2; Rec. 7)	Lack of evaluation (Rec. 1); Lack of understanding as to how to use PLA (Rec. 2; Rec. 7)	Lack of systematic evaluation (contradictory outcomes) (Rec. 1; Rec. 7)	Lack of evaluation (Rec. 1; Rec. 7)	Course design should define PLA use	Alignment across stakeholders (Rec. 3); Development of digital skills
Factors explaining slow uptake	RQ2	Teachers’ workload (Rec. 5); Varied course designs	Institutional changes impacting teachers’ work (Rec. 5); Lack of evidence about PLA effectiveness (Rec. 1)	Lack of ongoing support; Development of relevant skills	Management priorities (Rec. 4); Investment in staff	Lack of evidence (Rec. 1); Voluntary nature of participation; Lack of training; Teachers’ contracts (Rec. 5)	Lack of a common vision about PLA (Rec. 6)

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

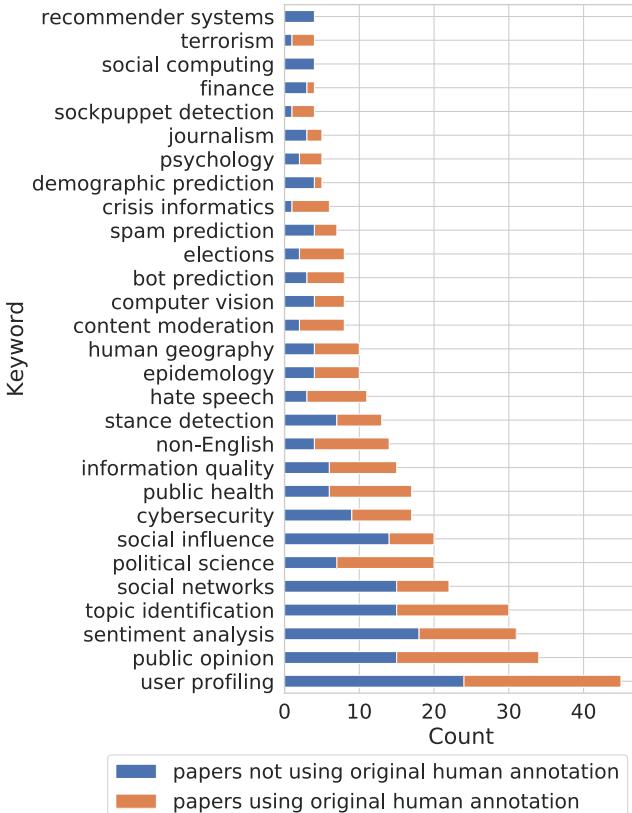
Possibilities

Conclusion  
and  
discussion

References

# Qualitative rigor?

- “Garbage in, garbage out?  
Do machine learning application papers in social computing report where human-labeled training data comes from?”
- Where labels come from matter!



# Qualitative rigor in context, too?

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

- "Contextual analysis of social media: The promise and challenge of eliciting context in social media posts with natural language processing"
- "VATAS: An Open-Source Web Platform for Visual and Textual Analysis of Social Media": "The lack of expertise is more relevant for marginalized communities or other critical domains that require specific training or background knowledge that are common in the field of social work. Poor performance is subsequently propagated into machine-learning models, as the models statistically fit the resulting data set with the purpose of learning to label samples the same way it was done by annotators. As a result, unreliable annotations can lead to models with low classification accuracy and biased predictions. This issue is why social work should drive social media annotation and interpretation of data and results, particularly when it relates to the most challenging social problems."

Labels	Loss			Other			Aggression			Macro F1
	p	r	f	p	r	f	p	r	f	
Gold	77.08	56.92	65.49	88.04	95.76	91.74	50	27.59	35.56	64.26
Distant	50.00	48.46	49.22	85.63	84.50	85.06	19.72	24.14	21.71	52.00

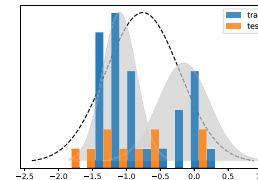
Table 2: SVM performance trained on hand-labeled vs distantly-labeled data.  
The difference between F1 scores is statistically significant with p=0.001.

[Overview](#)
[Background/  
review](#)
[Scope;  
review of  
some  
critiques](#)
[Examples of  
work that  
take, or  
enable  
critique](#)
[Possibilities](#)
[Conclusion  
and  
discussion](#)
[References](#)

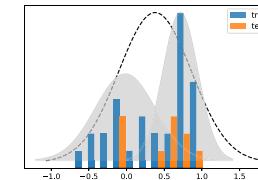
# Adding complexity?

- “Inherent disagreements in human textual inferences”
- We can always add complexity, but it is seldom a good idea for the purposes of modeling: when, and how, is it worth the extra bother?

p: A homeless man being observed by a man in business attire.  
h: Two men are sleeping in a hotel.



p: Paula swatted the fly.  
h: The swatting happened in a forceful manner.



p: Someone confessed that a particular thing happened.  
h: That thing happened.

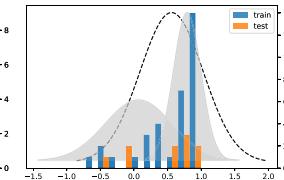


Figure 4: Examples of sentence pairs with bi-modal human judgment distributions. Examples are drawn from SNLI, the VerbNet portion of DNC, and the MegaVerdicality portion of DNC (from left to right). Training distribution is in blue; test in orange. Dotted black line shows the model fit when using a single component; shaded gray shows the model learned when allowed to fit  $k$  components. Distributions are over z-normalized scores in which 0 roughly corresponds to neutral ( $p \not\rightarrow h$ ) but not precisely ( 3.3).

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Investigating context?

- “Energy and Policy Considerations for Deep Learning in NLP”
- The paper that first pointed out the enormous energy cost (and therefore, CO<sub>2</sub> emissions) of deep learning models

<b>Consumption</b>	<b>CO<sub>2</sub>e (lbs)</b>
Air travel, 1 person, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000

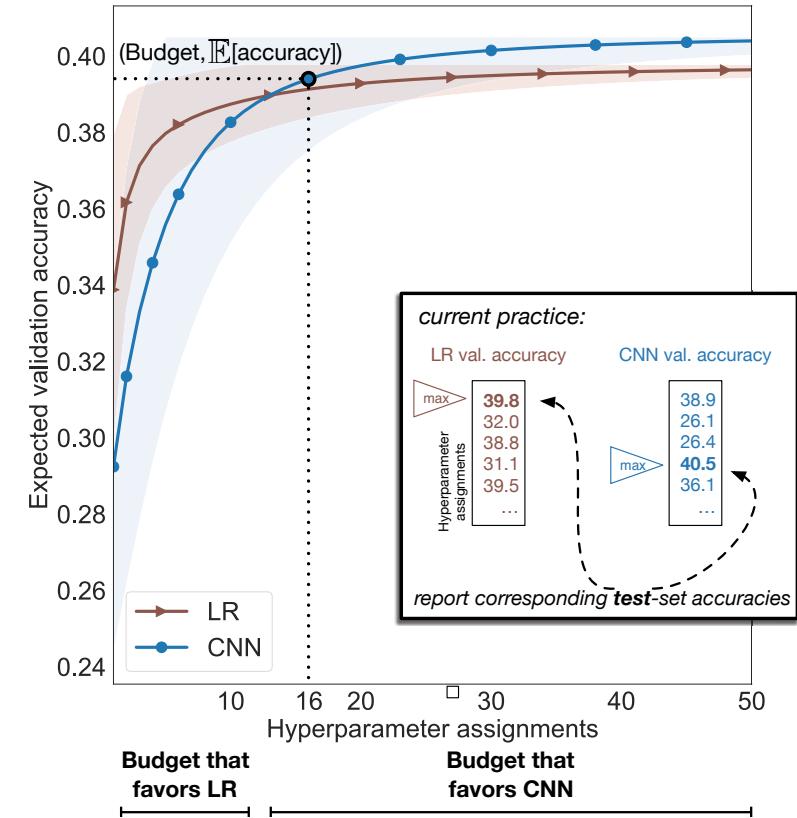
  

<b>Training one model (GPU)</b>	
NLP pipeline (parsing, SRL)	39
w/ tuning & experiments	78,468
Transformer (big)	192
w/ neural arch. search	626,155

Table 1: Estimated CO<sub>2</sub> emissions from training common NLP models, compared to familiar consumption.<sup>1</sup>

# Parameterizing energy/computation budget?

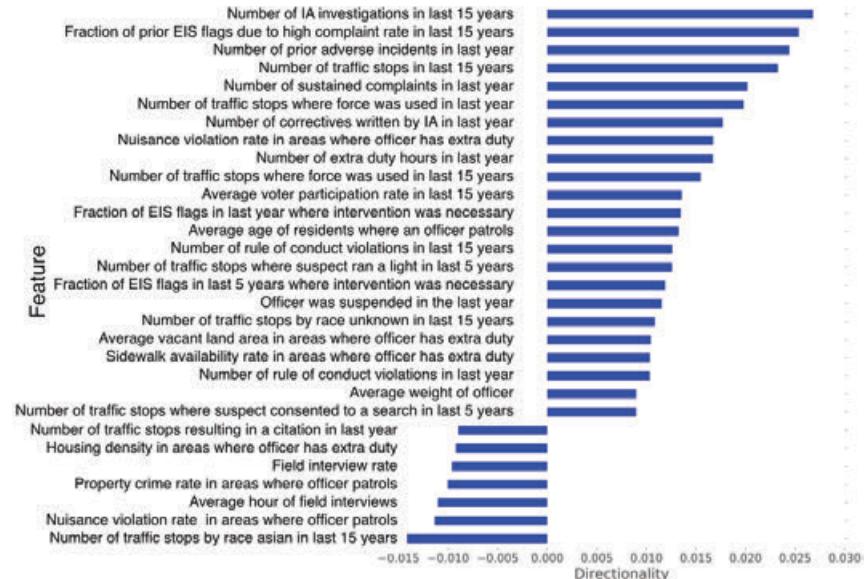
- “Show Your Work: Improved Reporting of Experimental Results”
- If you optimize for accuracy and energy usage, different models come out on top



[Overview](#)
[Background/  
review](#)
[Scope;  
review of  
some  
critiques](#)
[Examples of  
work that  
take, or  
enable  
critique](#)
[Possibilities](#)
[Conclusion  
and  
discussion](#)
[References](#)

# Studying up?

- “Identifying police officers at risk of adverse events”
- I.e., police brutality
- Why not subject those in power to the same “algorithmic” surveillance?
- (See also: “Studying up: reorienting the study of algorithmic fairness around issues of power”)



# Studying the whole system?

- “Algorithmic risk assessment in the hand of humans”
- Judges can ignore recommendations. They do so... but not at random
- Gives justification for giving harsher sentences to Black people (“the ‘algorithm’ says so!”)
- But judges exercise discretion with younger defendants (the innocence of youth)
- (Also looks at the aggregate impact of the system)

	(1)	(2)	(3)	(4)	(5)	(6)
	Panel A: Diverted   risk = low					
Alternative risk score	0.013 (0.010)					0.010 (0.010)
Black		-0.015 (0.015)				-0.014 (0.016)
Unemployed			0.025 (0.017)			0.009 (0.018)
Female				0.040** (0.016)		0.038** (0.017)
Age<23					0.069*** (0.020)	0.065*** (0.020)
Observations	3943	3943	3943	3943	3943	3943
R <sup>2</sup>	0.204	0.204	0.204	0.205	0.206	0.280
Mean DV	0.44	0.44	0.44	0.44	0.44	0.44
	Panel B: Diverted   risk = high					
Alternative risk score	-0.004 (0.005)					-0.007 (0.005)
Black		-0.029*** (0.010)				-0.045**** (0.012)
Unemployed			0.043**** (0.012)			0.018 (0.012)
Female				0.038*** (0.013)		0.040*** (0.014)
Age<23					0.065*** (0.011)	0.058**** (0.011)
Observations	7598	7598	7598	7598	7598	7598
R <sup>2</sup>	0.142	0.143	0.144	0.143	0.146	0.197
Mean DV	0.16	0.16	0.16	0.16	0.16	0.16

\* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01, \*\*\*\* p < 0.001

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Irony?

- “Predicting Financial Crime: Augmenting the Predictive Policing Arsenal”
- Not something that we’d expect will actually be taken up: but this is what it would look like if we treated white people’s crimes as symmetric with those of Black people

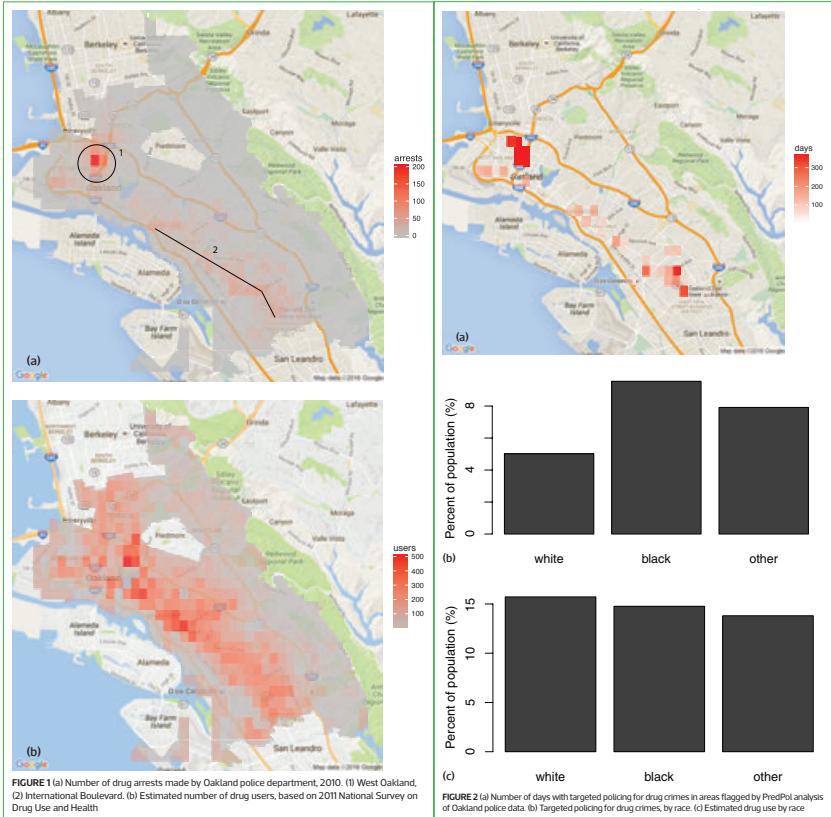


Fig. 6: The WCCEWS user interface. The map shows downtown Manhattan in New York City, NY. Color indicates the predicted density of white collar criminal activity. The left-hand panel shows the Top Risk Likelihoods of the listed crimes occurring within the selected geohash. Below is a histogram indicating predicted Approximate Crime Severity associated with discrete brackets of violation amount in \$USD. Finally, the panel lists Potential Offenders operating within the selected geohash, and a generalized white collar criminal subject.

[Overview](#)
[Background/  
review](#)
[Scope;  
review of  
some  
critiques](#)
[Examples of  
work that  
take, or  
enable  
critique](#)
[Possibilities](#)
[Conclusion  
and  
discussion](#)
[References](#)

# Re-application?

- “To predict and serve?”
- Re-applying PredPol on the basis of national surveys of drug use, rather than prior “crime” data



Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

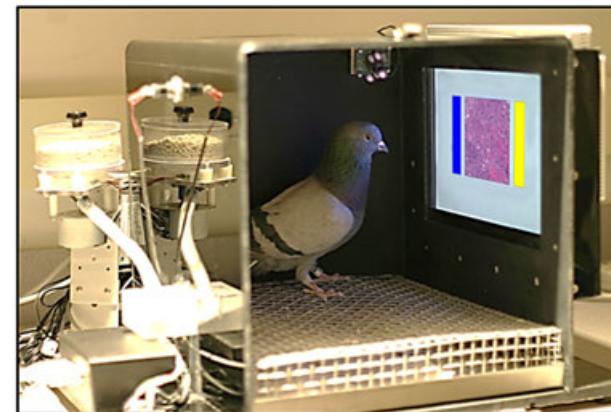
Possibilities

Conclusion  
and  
discussion

References

# Accidental absurdity? (*Or is it?*)

- “Pigeons as trainable observers of pathology and radiology breast cancer images”
- Pigeons competitive with deep learning (as another paper points out)



**Fig 1. The pigeons' training environment.** The operant conditioning chamber was equipped with a food pellet dispenser, and a touch-sensitive screen upon which the medical image (center) and choice buttons (blue and yellow rectangles) were presented.

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

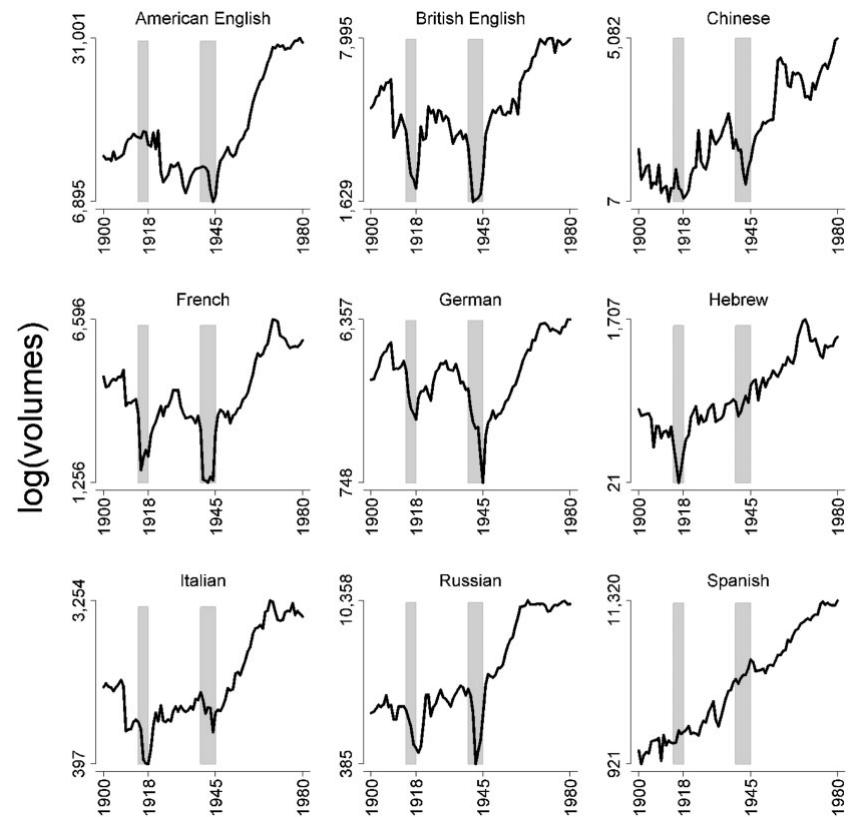
Possibilities

Conclusion  
and  
discussion

References

# Clean up?

- “The impact of lacking metadata for the measurement of cultural and linguistic changes using the Google Ngrams data sets—Reconstructing the composition of the German corpus in times of WWII”
- Critique of “culturomics” paper that uses Ngram frequencies to make sociological and historical claims
- (See also: “Characterizing the Google Books Corpus: Strong limits to inferences of socio-cultural and linguistic evolution”)



[Overview](#)
[Background/  
review](#)
[Scope;  
review of  
some  
critiques](#)
[Examples of  
work that  
take, or  
enable  
critique](#)
[Possibilities](#)
[Conclusion  
and  
discussion](#)
[References](#)

# Accidental reflexivity?

- “Human Rademacher Complexity”
- Had people try to remember/“learn” random labels: quantifies how we interpret signal in pure noise
- Clear about its assumptions: that everybody has the same “bounds”, but also that this should be relaxed

**Theorem 1.** Let  $\mathcal{F}$  be a set of functions mapping to  $[-1, 1]$ . For any integers  $n, m$ ,

$$P \left\{ \left| R(\mathcal{F}, \mathcal{X}, P_X, n) - \frac{1}{m} \sum_{j=1}^m \sup_{f \in \mathcal{F}} \left| \frac{2}{n} \sum_{i=1}^n \sigma_i^{(j)} f(x_i^{(j)}) \right| \right| \geq \epsilon \right\} \leq 2 \exp \left( - \frac{\epsilon^2 nm}{8} \right) \quad (2)$$

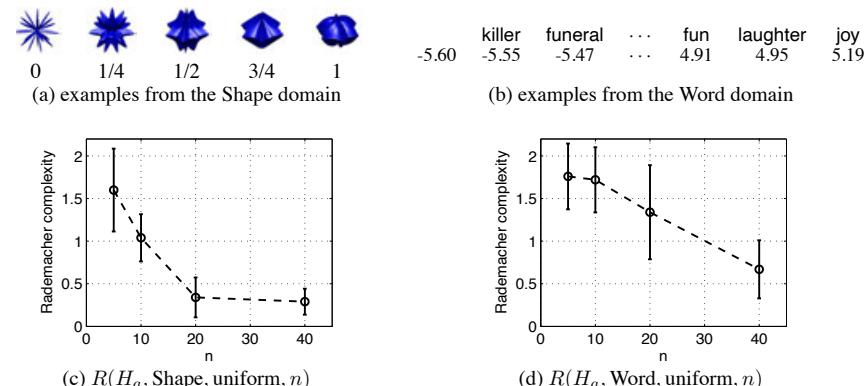


Figure 1: Human Rademacher complexity on the “Shape” and “Word” domains.

# Comparative study?

- “The Future of Coding: A Comparison of Hand-Coding and Three Types of Computer-Assisted Text Analysis Methods”

**Table 1.** Evaluation Metrics for Each Automated Method.

Method (Coding Scheme)	Inequality/Relevant			Not Inequality/Irrelevant			Weighted Average <sup>a</sup>			Time Trends		Support
	Precision (1)	Recall (2)	FI-Score (3)	Precision (4)	Recall (5)	FI-Score (6)	Precision <sup>b</sup> (7)	Recall <sup>b</sup> (8)	FI Score <sup>b</sup> (9)	(Two-Year MA) (10)	Correlation (11)	
(1) Supervised ML <sup>c</sup>												
Relevant versus irrelevant (A)	.85	.90	.87	.81	.74	.77	.83 (.81–.86)	.84 (.81–.86)	.83 (.81–.86)	.75	.74	745
Inequality versus not inequality (B)	.73	.60	.66	.80	.88	.84	.78 (.74–.80)	.78 (.75–.80)	.78 (.74–.80)	.69	.63	745
Inequality versus economic versus irrelevant (C)	.67	.70	.69	.76	.84	.80	.68 (.64–.71)	.69 (.65–.71)	.69 (.64–.71)	.72	.69	745
(2) Dictionary												
Levay-Enns (D)	.91	.25	.40	.83	.99	.91	.85	.84	.80	.42	.59	1,253
McCall (B)	.48	.84	.61	.86	.52	.65	.73	.63	.64	.66	.44	1,253
(3) Unsupervised ML												
Topic model versus explicit (D)	.63	.45	.53	.86	.93	.90	.82	.83	.82	.58	.68	1,253
k-means versus explicit (D)	.88	.14	.24	.81	.99	.89	.83	.81	.76	N/A	N/A	1,253

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Exploration?

- “Race, Writing, and Computation: Racial Difference and the US Novel, 1880-2000”
- “Computational methods demand the quantification of one’s objects of study. It’s likely easier to accept measuring a novel’s popularity by sales figures or classifying its genre by diction than labeling it according to discrete racial identifiers. Such labeling is an affront to critical race studies, which has taken as its very mission the deconstruction of racial categories.”
- “Critical suspicion, of course, can also lead to critical adaptation.”
- Carefully done, but underwhelming in the end

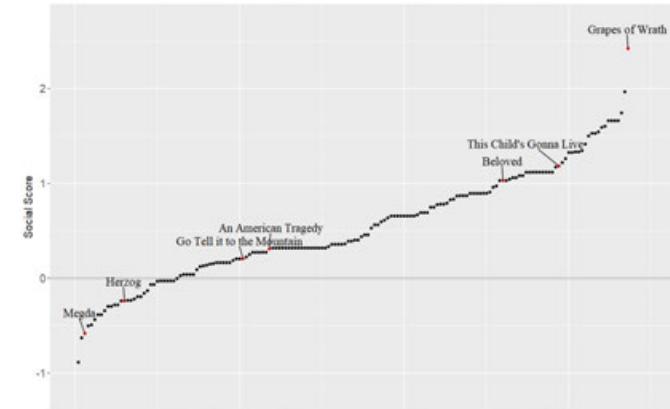
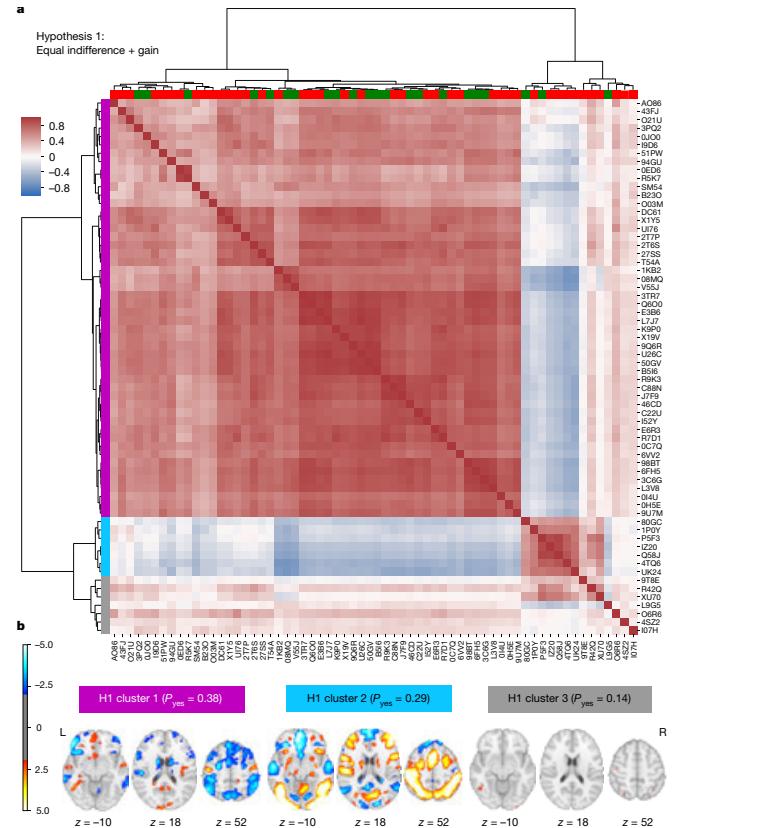


Figure 3. Plot showing the “social” score for all novels containing alignments with the Bible. Lower scores indicate novels where the Bible is less frequently cited in a “social” way, as we define the term. Scores closer to zero indicate novels where the “social” and “non-social” contexts are split evenly, as in James Baldwin’s Go Tell it to the Mountain.

[Overview](#)
[Background/review](#)
[Scope; review of some critiques](#)
[Examples of work that take, or enable critique](#)
[Possibilities](#)
[Conclusion and discussion](#)
[References](#)

# Investigating inconsistency?

- “Variability in the analysis of a single neuroimaging dataset by many teams”
- Lots of variation in analytic pipeline, and in results (but they suggest meta-analysis can resolve)
- (See also: “Many analysts, one data set: Making transparent how variations in analytic choices affect results”)



Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Using math where appropriate

- “Recombination: A family of Markov chains for redistricting”
- (Not actually “machine learning”, but actually algorithms, specifically MCMC)
- Compare proposed redistricting to distributions of random valid redistricting

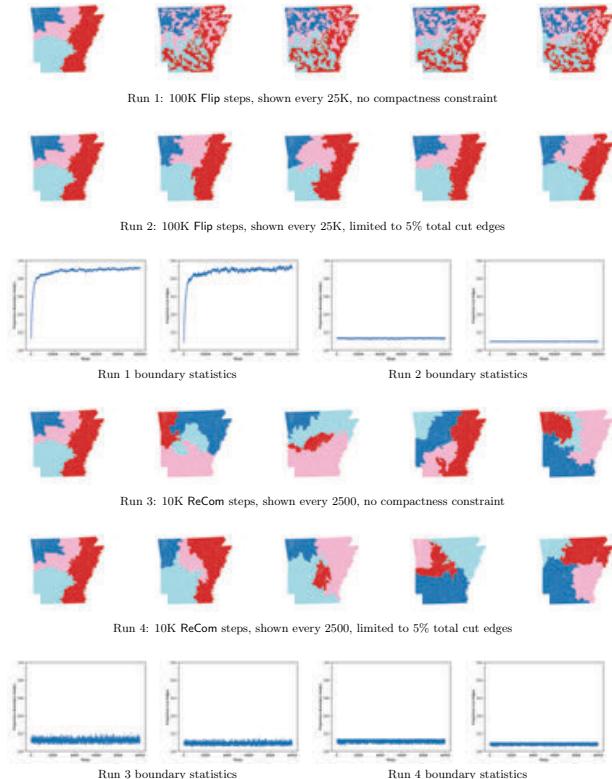


Figure 7: Arkansas block groups partitioned into  $k = 4$  districts, with population deviation limited to 5% from ideal. Imposing a compactness constraint makes the Flip chain unable to move very far.

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Building systems?

- “Alex speaks with my voice!”  
Promoting science discourse with bidialectal virtual peers”
- An AAVE virtual conversational agent



Overview

 Background/  
review

 Scope;  
review of  
some  
critiques

 Examples of  
work that  
take, or  
enable  
critique

Possibilities

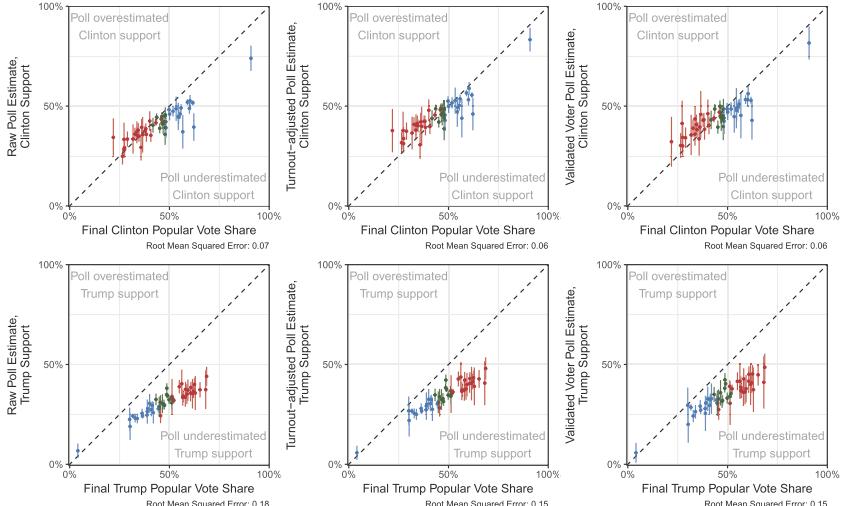
 Conclusion  
and  
discussion

References

# Statistical critique?

- “Statistical paradises and paradoxes in big data (I): Law of large populations, big data paradox, and the 2016 US Presidential Election”
- Derives a new fundamental statistical equation quantifying the effect of data quality/bias

$$\overline{G}_n - \overline{G}_N = \underbrace{\rho_{RG}}_{\text{Data Quality}} \times \underbrace{\sqrt{\frac{1-f}{f}}}_{\text{Data Quantity}} \times \underbrace{\sigma_G}_{\text{Problem Difficulty}}$$



**FIG. 4.** Comparison of actual vote shares with CCES estimates (and 95% confidence interval) across 50 states and DC. Top row for Clinton; bottom row for Trump. Color indicates a state's partisan leanings in 2016 election: solidly Democratic (blue), solidly Republican (red), or swing state (green). The left plot uses sample averages of the raw data ( $n = 64,600$ ) as estimates; the middle plot uses estimates weighted to likely voters according to turnout intent (estimated turnout  $n = 48,106$ ); and the right plot uses sample averages among the subsample of validated voters (subsample size  $n = 35,829$ ). Confidence intervals based on unweighted sample proportions are computed following (3.9), where the use of SRS variances can be conservative given the stratified design of the survey, and yet they still do not provide any realistic protection against the increased MSE caused by the non-response bias. For the turnout adjusted estimate, which is in a ratio form, a  $\delta$ -method is employed to approximate its variance, which is then used to construct confidence intervals.

[Overview](#)
[Background/  
review](#)
[Scope;  
review of  
some  
critiques](#)
[Examples of  
work that  
take, or  
enable  
critique](#)
[Possibilities](#)
[Conclusion  
and  
discussion](#)
[References](#)

# Statistical critique?

- “The grand leap”
- (See also:  
“Graphical models  
for causation, and  
the identification  
problem”)

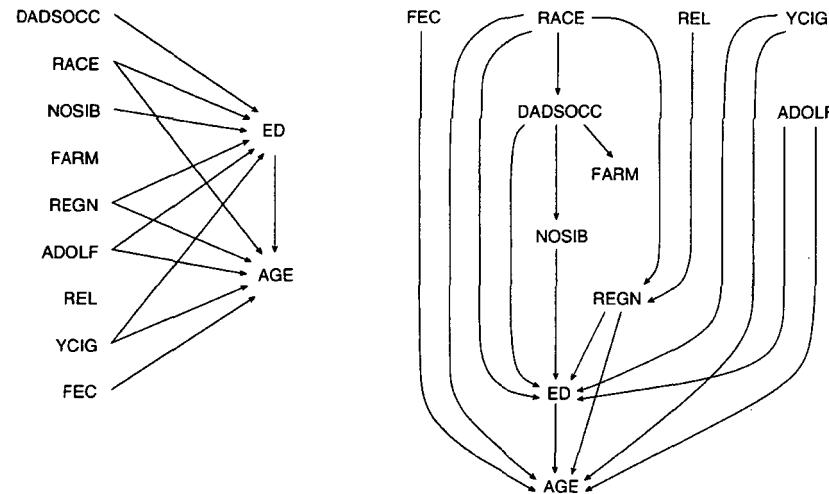


Fig.3. The left-hand panel shows the model reported by SGS (p.140). The right-hand panel shows the whole graph produced by the SGS search program TETRAD.<sup>16</sup>

# Statistical critique?

Overview

Background/  
review

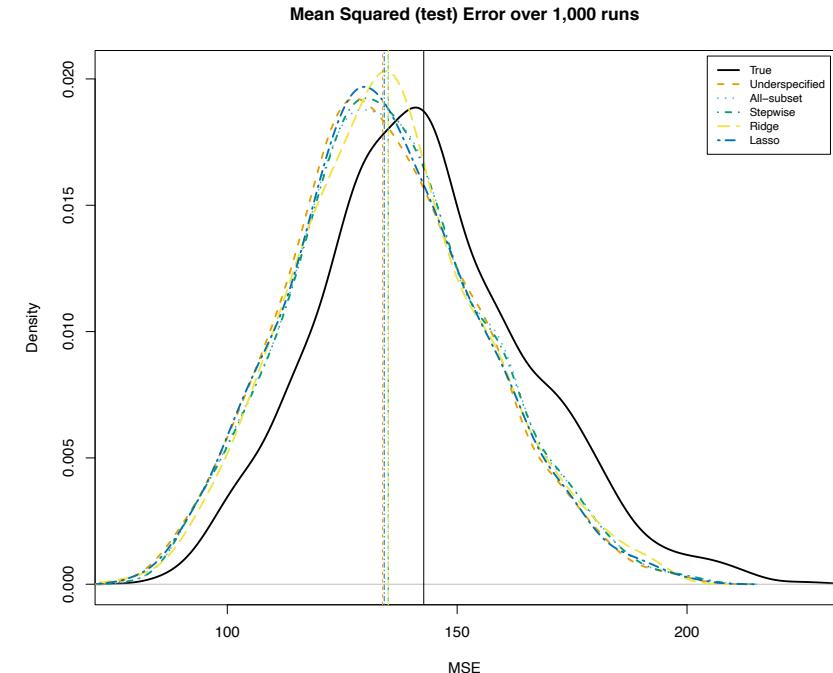
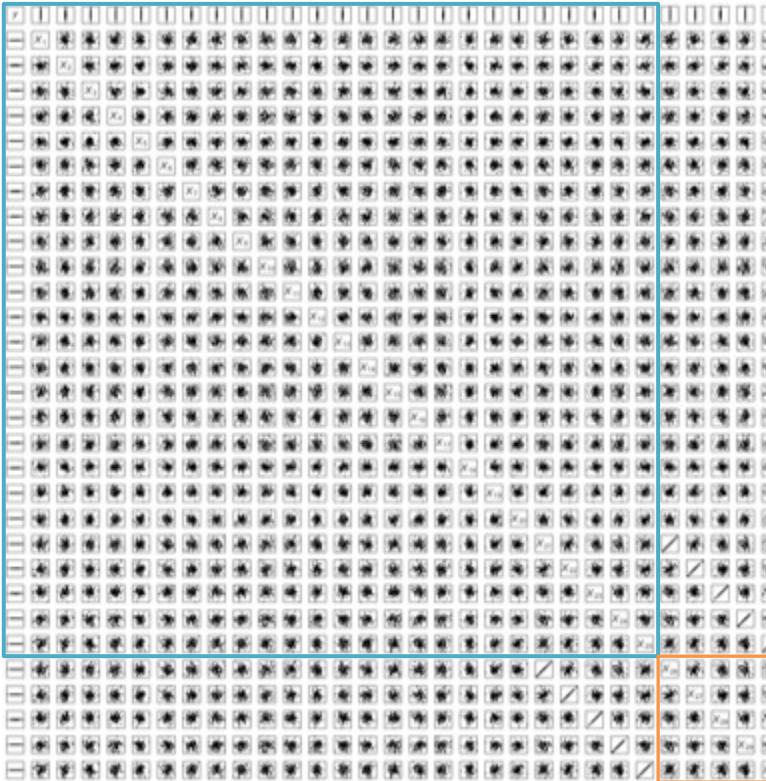
Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References



Simulation of Shmueli, 2010, *Stat. Sci.*

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

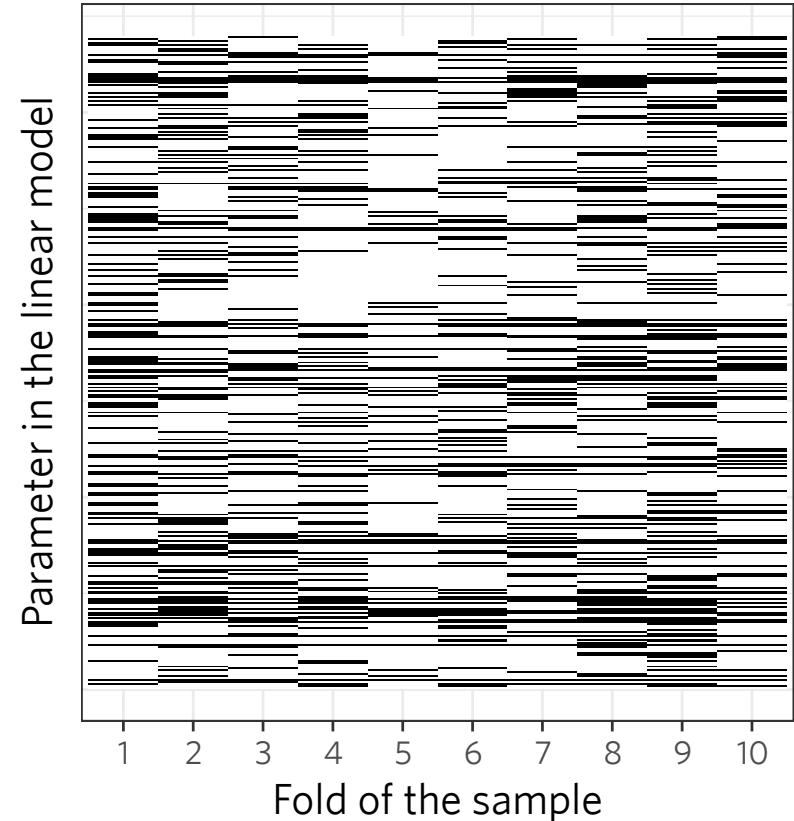
Possibilities

Conclusion  
and  
discussion

References

# Statistical clarification?

- “Machine learning: An applied econometric approach”
- Very different sets of correlations can “predict” (fit) equally well
  - Leo Breiman (2001) called this the “Rashomon Effect”
- But different fits suggest very different interventions



Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

## **Possibilities**

Conclusion  
and  
discussion

References

# **Possibilities**

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

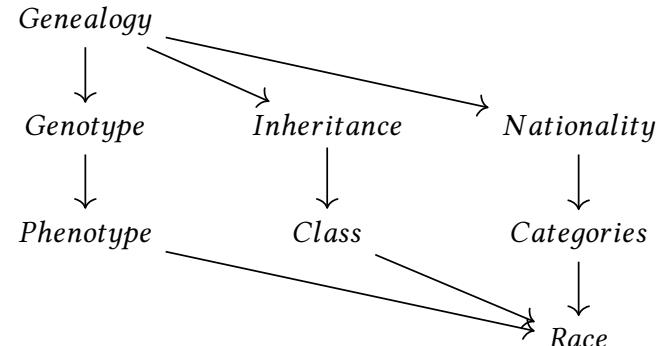
Possibilities

Conclusion  
and  
discussion

References

# Graphical models to express

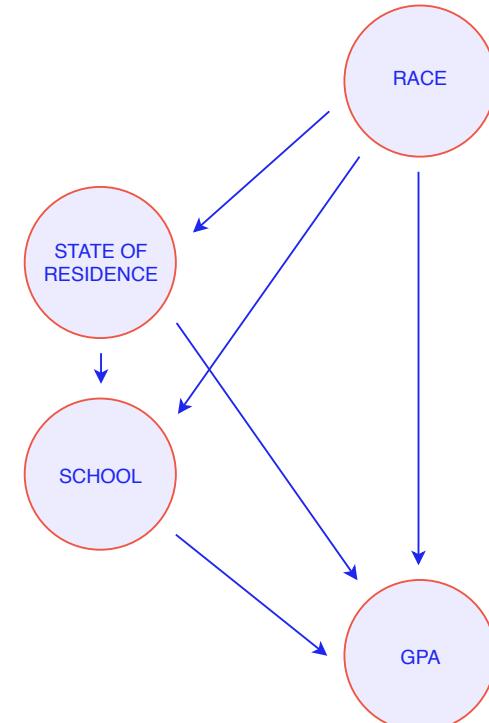
- “Racial categories in machine learning”
- (See also: “Towards a critical race methodology in algorithmic fairness”)



**Figure 1:** A model of how individual biological properties (genealogy, genotype, and phenotype) are racialized through national political categories and associations with socioeconomic class. Here inheritance refers to all forms of capital, including economic and social, passed from parents to children. Broadly speaking, genealogy is a strong determiner of race, but importantly as a common cause of phenotype, class, and nationally recognized racial categories, which are separate components of racial classification.

# Graphical models as guide to critique

- “Disparate causes, pt. I” and “Disparate causes, pt. II”
  - Would need to backtrack through entire life history to really account for race as a counterfactual
- (See also: “Eddie Murphy and the dangers of counterfactual reasoning”)



Hu 2019b

# Reparations and COVID-19 paper

Overview

Background/  
review

Scope;  
review of  
some  
critiques

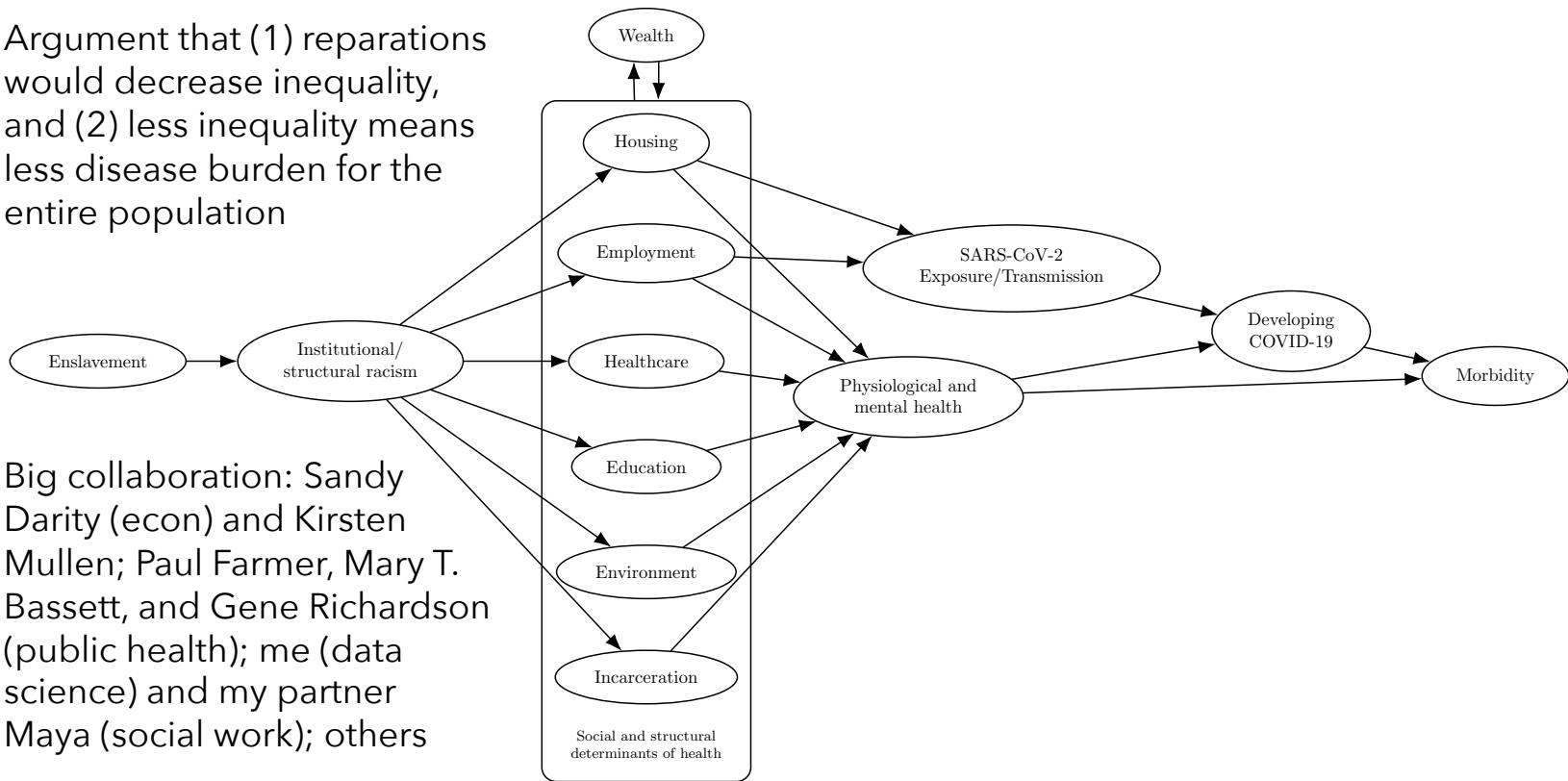
Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

Argument that (1) reparations would decrease inequality, and (2) less inequality means less disease burden for the entire population



# Classic argument for CV

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

**Possibilities**

Conclusion  
and  
discussion

References

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

Irreducible      Bias squared      Variance      “Optimism”

Overview

Background/  
reviewScope;  
review of  
some  
critiquesExamples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

# Apply this to non-iid data

- 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

Overview

Background/  
review

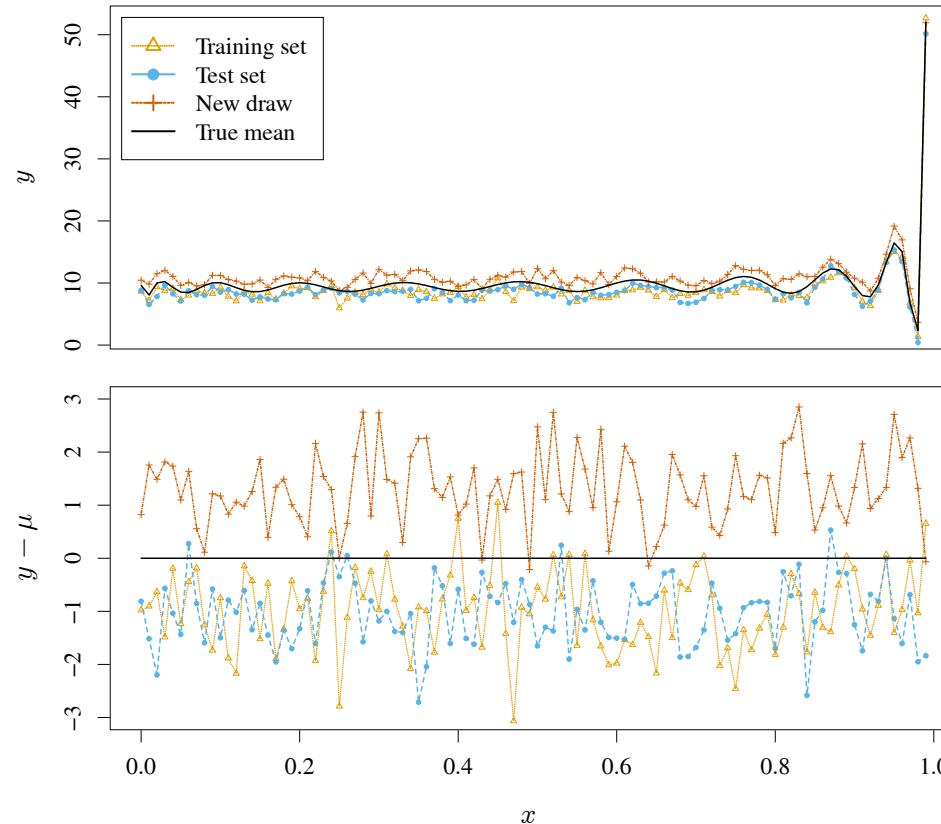
Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

## Possibilities

Conclusion  
and  
discussion

References



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

Overview  
Background/  
review

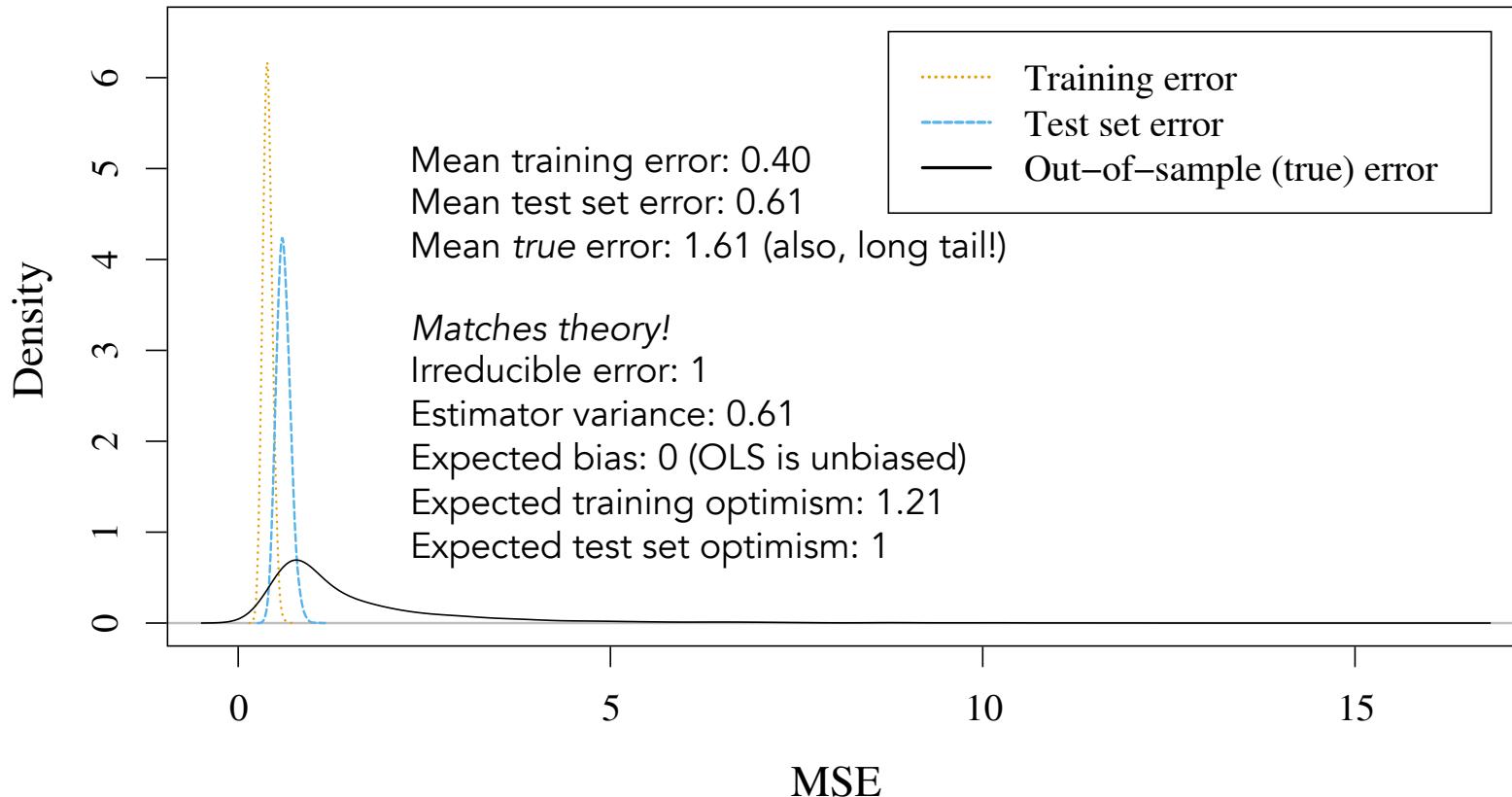
Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

## Possibilities

Conclusion  
and  
discussion

References



Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

**Conclusion  
and  
discussion**

References

# Conclusion and discussion



# Some things we covered

- Auditing
- Scalable measurement
- Qualitative follow-up
- Qualitative evaluation
- Qualitative rigor (in context)
- Adding appropriate complexity
- Investigating context
- Parameterizing budget
- Studying up
- Studying the whole system
- Re-application
- Absurdity
- Clean-up
- Reflexivity
- Comparative study
- Exploration
- Investigating inconsistency
- Using math where appropriate
- Building systems
- Statistical critique/clarification

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References



# What's the point?

Overview

Background/  
review

Scope;  
review of  
some  
critiques

Examples of  
work that  
take, or  
enable  
critique

Possibilities

Conclusion  
and  
discussion

References

With increasing depth:

- Rigor. Just doing things more carefully and outlining limitations (including sociotechnical)
- Critiquing/destabilizing machine learning: listing limitations, maybe absurdity (!)
- Strategic positivism: using machine learning to add legitimacy (risks reifying)
  - (See also: "Roles for computing in social change": but maybe not strategic, maybe just positivism)
- Imagining, mixed methods