

Artificial Intelligence/Operations Research Workshop 2 Report Out

Workshop

August 16-17, 2022

Workshop Organized by

John Dickerson, University of Maryland
Bistra Dilkina, University of Southern California
Yu Ding, Texas A&M University
Swati Gupta, Georgia Institute of Technology
Pascal Van Hentenryck, Georgia Institute of Technology
Sven Koenig, University of Southern California
Ramayya Krishnan, Carnegie Mellon University
Radhika Kulkarni, SAS Institute, Inc. (retired)

With Support From

Catherine Gill, Computing Community Consortium (CCC)
Haley Griffin, Computing Community Consortium (CCC)
Maddy Hunter, Computing Community Consortium (CCC)
Ann Schwartz, Computing Community Consortium (CCC)



CCC

Computing Community Consortium
Catalyst



Background and Introduction

Artificial intelligence (AI) has received significant attention in recent years, primarily due to breakthroughs in game playing, computer vision, and natural language processing that captured the imagination of the scientific community and the public at large. Many businesses, industries, and academic disciplines are now contemplating the application of AI to their own challenges. The federal government in the US and other countries have also invested significantly in advancing AI research and created funding initiatives and programs to promote greater collaboration across multiple communities. Some of the investment examples in the US include the establishment of the [National AI Initiative Office](#), the launch of the [National AI Research Resource Task Force](#), and more recently, the establishment of the [National AI Advisory Committee](#).

In 2021 INFORMS and ACM SIGAI joined together with the Computing Community Consortium (CCC) to organize a series of three workshops. The objective for this workshop series is to explore ways to exploit the synergies of the AI and Operations Research (OR) communities to transform decision making. The aim of the workshops is to establish a joint strategic research vision for AI/OR that will maximize the societal impact of AI and OR in a world being transformed by technological change and a heightened desire to tackle important societal challenges such as growing racial and social inequity, climate change, and sustainable solutions to our food-water and energy needs. The vision for the workshops is to exploit and expand on the emerging synergies between these two communities with complementary strengths. However, there are barriers and difficulties in realizing this vision due to cultural differences between AI and OR communities. The workshop series aims to overcome these difficulties and to provide a stepping stone for a strong and sustained collaboration between the two fields.

The first workshop was held virtually in September 2021 with speakers and participants drawn from leading researchers in both the Operations Research and Computing research communities. The expectations were to promote greater inter-disciplinary collaborations between the two areas, inspire the agenda for research, and address critical questions for the future of AI. The agenda included interspersed sessions on Methods and Applications with moderated Q&A and breakout group discussions. Some of the key outcomes of the workshop included:

- Dispelling common misconceptions across the two communities and discovering common objectives, albeit a few disparate assumptions leading to different tools and methods
- Exploring multiple opportunities for collaboration at the Data, Methods, and Policy layers

- Discussing the need for a science of making good decisions with collaboration across multiple disciplines in addition to AI and OR
- Listing some potential future directions including access to data sets for a variety of application domains, competitions requiring joint teams of OR and AI participants, starting an AI/OR summer school for students in both disciplines, etc.

For details including the presentations and summaries of the group discussions, see the [AI/OR Workshop I report](#).

Overall Theme for the Second Workshop

One key outcome from the first workshop was the decision that we need a deeper dive into the topics of Fairness and Ethics in AI and also include a discussion of the role of Causality in applications of AI as well as considerations of human computer interaction in the implementation of automated solutions. As a result, the second workshop was planned around Trustworthy AI with four sessions on related topics within the scope of AI and OR:

- Fairness
- Explainable AI / Causality
- Robustness / Privacy
- Human Alignment and Human-Computer Interaction

While the first workshop focused on articulating a strategic vision, this workshop focused on what it takes to develop and deploy intelligent systems based on AI and OR technology that are trustworthy. In particular, the workshop was designed to:

- Study the state of the art in Trustworthy AI from a multi-disciplinary lens and for various application domains.
- Articulate grand challenges that need to be overcome to deploy trustworthy AI systems in the wild. Specifically, what tools and technologies have to be developed and evaluated for each of the foundational elements of trustworthy AI systems.
- Select a few topics for summer schools and research programs to foster collaborations in AI/OR for these topics.

The speakers and participants were selected from computer science and OR communities to foster a healthy exchange of ideas between the two groups. A key requirement for this workshop was to ensure that all speakers as well as participants attended it in person, to ensure a robust and healthy discussion on the topics presented

during the event. It was held in Atlanta, GA on the Georgia Tech campus on August 16-17, 2022.

This report out should be considered an interim report with a final report to be published after the end of the third workshop. We expect the next workshop to be held in the first half of 2023. Note that the workshops are held on behalf of the community, by invitation only.

The workshop was kicked off with welcoming remarks by Radhika Kulkarni and introductions of all the participants, followed by a few remarks from Mr. Murat Omay from the Department of Transportation on funding opportunities for OR and AI research, from a problem-driven perspective. The details of Mr. Omay's remarks and all the sessions are described in the remainder of the report.

Brief Comments from the Department of Transportation

Murat Omay - [Figure B-1: Overview of ITS JPO Programs](#)

After introductions, Murat Omay, from the Department of Transportation, gave a brief presentation on the DoT's Joint Program Office (JPO) for Intelligent Transportation Systems (ITS). Omay began by discussing the Data Access and Exchanges Portfolio, created and managed by the JPO ITS program. The portfolio aims to effectively generate, acquire, govern, manage and analyze ITS data across all modes to advance multimodal research and to enable a safe, equitable, multimodal and resilient transportation network. One of the key goals of this effort is to position the DoT in a leading role in data and AI strategies, data-based automation and AI research, innovation, and transformative data culture. Omay listed many of the current activities in this area, and explained how each furthered progress toward the overall goal of the JPO ITS program. Many of these activities were cross-departmental efforts, with ongoing dialogues between the DoT and the Office of Science and Technology Policy (OSTP), the Networking and Information Technology Research and Development program (NITRD) and the Open Knowledge Network (OKN) to name a few. Omay invited the participants to follow up with him directly to learn more about the Data Access and Exchanges Program.

Omay then turned our attention to the AI for ITS program, also run by the DoT. This program advocates for utilizing ethical AI and ML technologies to create a more efficient logistics system to transport goods and people. The program aims to coordinate technology and policy research to more quickly integrate AI and ML into the existing transportation system. In 2021, the AI for ITS program sought input from AI researchers on the "deployment ready" AI technologies the program had created, as well as on the best AI

areas for the DoT to invest in and any current DoT technologies which could utilize AI to improve their performance.

Omay listed some of the key challenges to AI adoption and implementation, which are not unique to ITS. These include explainability, liability, model drift, privacy, security, ethics and equity and others; addressing these challenges is an ongoing exercise. He called out the fact that maintaining a human-in-the-loop approach is helpful in identifying and mitigating these challenges. He acknowledged the need for a multi-pronged approach and welcomed input from the community to help in addressing these issues.

Omay then took questions from the workshop participants and asked for any feedback pertaining to the DoT's efforts to incorporate AI in their technologies. He recommended Robert Sheehan as the key point of contact to learn more about the AI for ITS program. Robert Sheehan (Robert.Sheehan@dot.gov) is Acting Chief of Policy, Architecture, and Knowledge Transfer, U.S. DOT ITS Joint Program Office.

Panel A: Fairness

Moderated by John Dickerson

Slide Presentations

Dimitris Bertsimas - [Figure B-2: Improving on Fairness/Bias](#)

Maria De-Arteaga - [Figure B-3: Social Norms Bias: Residuals Harms of Fairness-Aware Algorithms](#)

Nikhil Garg - [Figure B-4: Auditing and Designing for Equity in Government Service Allocation](#)

David Shmoys - [Figure B-5: Fairness as the Objective in Congressional Districting](#)

This panel broadly covered topics including fairness in allocation, learning, and decision making. Speakers broadly touched on the practical implementations of different definitions of fairness in, for example, healthcare, academic talent sourcing settings, and government resource allocation such as inspections of infrastructure. Speakers also discussed situations where traditional group fairness interventions do not work due to (by definition) averaging across an entire group. Fairness in redistricting to combat gerrymandering, a topic at the intersection of AI and OR research with economic and policy considerations as well, was also presented during this panel. A lively discussion followed considering, amongst other topics, the efficacy and appropriateness of using fairness as a *constraint* in various settings, and new topics in allocation and market design. Specifically, discussions of fairness in two-sided matching platforms such as rideshare (e.g., Uber and Lyft, where

drivers on one side of the market are matched to riders on the other) arose naturally, although that specific use case had not yet been covered in any of the talks.

We now dive deeper into each of the four talks. Each talk came from folks more on the “OR side” of the research world, although all four speakers have published in AI and/or more broadly-defined “AI/ML-first” venues such as the ACM Conference on Fairness, Accountability, and Transparency (FAccT), and each speaker maintains strong collaborations with those on the “CS side.”

The first talk, by Dimitris Bertsimas of MIT, primarily discussed an ongoing collaboration between MIT and Massachusetts General Hospital (MGH) analyzing a part of the healthcare pipeline, the discharge of trauma patients to post-acute care (PAC). Here, a group-level disparity was identified: 60% of trauma patients at MGH are Black, but only 12% of those receiving post-acute care are Black. Yet, it has been shown that PAC results in lower readmission rates and other benefits; this implies a potential disparate impact issue across racial categories. The talk covered in depth many of the intricacies of predicting risk and potential outcomes, and then presented a mixed-integer-programming-based (MIP-based) method for enforcing a form of demographic parity in PAC assignment. This approach also had interesting overlap with improving predictive performance of the assignment as well. The talk also connected the general approach used for this specific healthcare setting to, potentially, other settings such as talent sourcing in academic departments.

Perhaps most related to that first talk was Nikhil Garg of Cornell Tech’s third talk in the panel, where he gave a general framework for auditing government processes and subsequently working with various government agencies to improve on those processes. Roughly, the pitch was to Audit agency decisions along the entire pipeline, from problem description to data collection to data maintenance to model building to decisioning. Issues of efficiency and equitability arise throughout the full pipeline, and do not exist in isolation; that is, upstream deployment decisions feed into the current module, and that module’s outputs then impact downstream sections of the pipeline. As OR and AI practitioners, a particular focus here was on the challenge of modeling capacity-constrained decisions under various forms of uncertainty. The talk then went in-depth to an ongoing collaboration between Cornell Tech and New York City, looking to improve its 311-based method to crowdsource infrastructure incidents that require intervention from the city. Here, concerns of equity arise due to many forms of uncertainty (e.g., different areas of the city may have different reporting rates) and also allocative decisioning (e.g., areas close to city services may receive more effective treatment than those further away). A key takeaway from this talk was that focusing on fairness means comparing decisions for comparable incidents/groups – AI/ML tends to focus on analyzing individual incidents, whereas OR

tends to focus on making global decisions across groups, so a responsible approach to fairness necessarily requires techniques from both fields' core competencies.

The chronologically second talk in the panel was given by Maria de-Arteaga of UT Austin, which provided a nuanced and in-depth discussion of some surprising side effects of blindly applying traditional fairness interventions. Specifically, the talk focused on identifying residual harms due to biased inputs when running fairness-aware algorithms, such as those that may impose a group fairness concern as an explicit constraint in an optimization problem. The example given was in the natural language processing (NLP) space, where group membership in a study was inferred by automatically “reading” written biographies of various academics. That inference process was shown to be biased on its own, with those who identified as a particular gender *and* who followed societal norms regarding writing were more accurately identified as their correct group than those who did not. Then, using a fairness-aware algorithm that enforces group fairness of some type – but across those biased group labels that have been inferred by some upstream process – may in fact disproportionately and systematically harm specific types of individuals within their true groups. A key takeaway, then, is that the risks of individual harm can certainly exist even under strict adherence to group fairness definitions and constraints.

Finally, David Schmoys of Cornell closed out the session with a discussion of his group's recent work on gerrymandering and fairness in redistricting. Gerrymandering, at a high level, is a method of “unnaturally” shaping political districts such that the resulting vote aggregated across those districts artificially over- or under-represents a particular party. Fair redistricting, then, seeks to combat this form of unfairness by way of choosing a partitioning of a geospatial region that holds to some standards. This can be seen, and is often modeled, as a set partitioning problem in its ideal form, with a variety of geographic/geospatial desiderata such as compactness and demographic or representational desiderata such as proportional representation. This concluding talk discussed methods for defining and also implementing forms of fairness in redistricting, and kicked off a vibrant post-panel discussion amongst participants.

Networking Lunch

Brief remarks from Pascal Van Hentenryck (AI4OPT) and Andrew Kahng (TILOS)

During lunch Prof. Pascal Van Hentenryck and Prof. Andrew Kahng gave the participants a brief overview of their respective [AI Institutes](#), which have both received \$20 Million dollar grants from NSF for a period of five years. Below are brief descriptions of these two

institutes, [AI4OPT](#) and [TILOS](#), obtained from their websites and augmented with some comments from Pascal and Andrew.

TILOS

The TILOS mission is to make impossible optimizations possible, at scale and in practice. The institute's research will pioneer learning-enabled optimizations that transform chip design, robotics, networks, and other use domains that are vital to our nation's health, prosperity and welfare. TILOS is a partnership of faculty from the University of California, San Diego, Massachusetts Institute of Technology, National University, University of Pennsylvania, University of Texas at Austin, and Yale University. Many faculty members associated with TILOS are working on Fairness and Explainable AI, two of the themes of this workshop. TILOS is partially supported by the Intel Corporation.

AI4OPT

This NSF Artificial Intelligence (AI) Research Institute for Advances in Optimization aims at delivering a paradigm shift in automated decision-making at massive scales by fusing AI and Mathematical Optimization (MO), to achieve breakthroughs that neither field can achieve independently. The Institute is driven by societal challenges in energy, logistics and supply chains, resilience and sustainability, and circuit design and control. Moreover, to address the widening gap in job opportunities, the Institute delivers an innovative longitudinal education and workforce development program with an initial focus on historically black high schools and colleges in Georgia, as well as Hispanic-serving high-schools and colleges in California. The Institute is also developing internship programs with national laboratories and industrial partners, and is building a strong, welcoming, and inclusive community, highlighting social mobility opportunities and the societal impact of AI technologies. The focus is on fundamental and use-inspired research, with a goal of fusing machine learning and optimization by merging both the data-driven and model-based paradigms which are the hallmarks of the two areas of research.

Both of these institutes give junior and mid-level faculty the opportunity to lead by exploring the potential and challenges of use-inspired research in high-stakes domains. They also have strong collaborations with several industry partners listed on their websites.

Panel B: Human Alignment/HCXAI/HCI

Moderated by Swati Gupta

Slide Presentations

Hamsa Bastani - [Figure B-6: Decision-Aware Reinforcement Learning](#)

Peter Frazier - [Figure B-7: Preference Learning for Stakeholder Management](#)

Kristian Lum - [Figure B-8: De-biasing “Bias” Measurement](#)

Mark Riedl - [Figure B-9: Toward Human-Centered Explainable Artificial Intelligence](#)

Panel B broadly consisted of decision making in various contexts which need to account for human factors, such as interpretability, ways in which models are explained, adaptability of the models, and changes from the status quo.

The panel started with Hamsa Bastani (Wharton) presenting her work with Sierra Leone National Medical Supplies Agency, where they focused on the distribution of essential medicines and managing inventory at different health facilities in the region, while navigating highly uncertain demands. Hamsa talked about aligning the loss function used to train the machine learning model with the decision loss associated with the downstream optimization problem. She interpreted the gradient of their loss function as a simple re-weighting of the training data, allowing it to flexibly and scalably be incorporated into complex modern data science pipelines, yet producing sizable efficiency gains. She shared results showing the decrease in unmet demand of essential medicines with the use of a decision-aware loss function versus a decision-blind loss function.

Next, Peter Frazier (Cornell and Uber) discussed understanding stakeholders and their constraints, in the context of COVID testing policies for Cornell University. Peter posed the question: Can AI help with tasks that OR practitioners need to do manually? For example, can AI help us understand stakeholder preferences? He advocated that Bayesian Optimization is a black-box derivative-free non-convex optimization method, which can be used for preference learning to get utility functions from stakeholders. The perspective of using AI to even understand the stakeholder preferences (which they may not be able to even quantify as functions or utilities themselves) was perceived as an interesting human-AI use case that has not received enough attention. The Cornell study showed that there is plenty of room for AI-enabled stakeholder engagement for OR applications. In particular, this approach helps one:

- Understand stakeholder goals, beliefs and incentives
- Understand how groups of stakeholders influence each other
- Predict how stakeholders will react to communication and
- Manage trust (in the OR analyst and her models).

The third talk in the session was given by Kristian Lum (Twitter), where she discussed the implications of checking for bias in terms of any algorithm's impact on different demographic groups. She postulated that while much of the work in algorithmic fairness over the last several years has focused on developing various definitions of model fairness (the absence of group-wise model performance disparities) and eliminating such "bias," much less work has gone into rigorously measuring it. She argued that many of the metrics used to measure group-wise model performance disparities are themselves statistically biased estimators of the underlying quantities they purport to represent. For example, the amount of "bias" measured increases as the number of groups increase, but this is just statistical noise. Lum proposed a "double-corrected" variance estimator to use instead. The key message of her talk was to question the statistical properties of various metrics for fairness used in practice and research, and that a lot of attention is needed to understand what we mean by fairness, inequity, and diversity. Meta-metrics cannot capture the entirety of the impact of ML systems. Small measured disparities should not be taken as a guarantee that the system is fair or free from adverse impacts.

Finally, we ended the session with a thoughtful talk by Mark Reidl (Georgia Tech), on explainable AI for consumer facing AI, and how different ways of communicating explanations themselves have a different impact on the users. For example, non-experts value contextual accuracy, awareness, strategic detail, intelligibility, relatability, and had unwarranted faith in numbers (they would trust the system more when numbers were used in the explanation). He claimed that our goal should not be to get people to trust the AI but to appropriately trust the AI via trust calibration. He discussed new research in explainable reinforcement learning and experiential explanations.

Overall the session touched upon various aspects of human-AI collaboration: how to interpret models for decision-makers, how to understand stakeholder preferences, how to even understand bias predictors, as well as the importance of designing good explanations to increase trust in AI.

Breakouts for Panels A and B

Breakout 1

Moderator: Jon Owen

Notetaker: Catherine Gill

- **Human Interactions:** We discussed an interesting opportunity of expanding the conversation beyond AI and OR communities to include disciplines related to

human decision making and psychology. Several motivating examples were considered, including: (a) the risk of introducing bias through the questions we ask and how we ask them when forming our understanding of a system (e.g., OR modeling) or executing market research that guides modeling and provides input data, (b) the evaluation and verification of results and outcomes, and (c) our effectiveness towards influencing actions through AI/OR analyses and the communication of results and implications – including human-in-the-loop decision making. It was mentioned that there are 3 places where humans come in: early, like choosing/inputting data; in the middle, where decisions get made; and the end, in evaluation and verification. At the end of the day, we want recommendations to be robust and to make sense to the decision maker.

- **Public Perceptions of AI:** We also discussed public misperceptions about AI due to hype. It was noted that many naïve users view AI as almost “magical”, and that it will always be better at solving problems than humans. This led to discussion around (a) the need for a better understanding of these techniques as tools, not solutions, and (b) greater transparency around embedded assumptions and their potential implications on interpretation of results, including the role of bias.
- **Data Quality:** We discussed several points related to data and notions of data quality. Often you don’t know the quality of data, especially when underlying data generating functions are poorly understood, or a specific use case for the data differs from the original intent (e.g., when data is generated/collected for one purpose but used for another). With these limitations, the need for data learning was discussed, as well as an observed gap between recording state actions without recording the associated probabilities.
- **AI “Intelligence” and Automation:** The concept of AI “intelligence” was briefly discussed as both an aspiration and a moving target based on our expectations for automation. Should we care about automating everything through AI/ML? Also, there are limits for data-centric approaches moving beyond descriptive and predictive; for example, prescriptive recommendations require going beyond the data to be useful – this is an opportunity for AI and OR to combine.

Breakout 2

Moderator: Harrison Schramm

Notetaker: Haley Griffin

- **Fairness and Bias:** We discussed fairness and the many ways in which bias can be introduced to a program. By attempting to limit the bias in a program a researcher may actually introduce more bias. Trying to limit unwanted confounding variables may result in detrimental unforeseen consequences, which can actually increase the unfairness of a program. Also, by reducing the bias of a program for

the majority of users, you may disproportionately increase the unfairness for minority groups of users, who may be the very people who need fairness protections the most.

- **Explainable AI:** We talked about explainable AI, and the need for explainable OR as well. Some participants took issue with the current state of explainable AI, insisting that it is overly vague, and the outputs of AI programs give users no information as to how to interpret these outputs. These participants advised creating more programs to analyze and interpret the results of AI programs to better inform users how an AI program came to a certain conclusion, rather than just returning an output with no explanation.

Breakout 3

Moderator: Theodora Chaspari

Notetaker: Ann Schwartz

- **Fairness:** Efficiency and fairness are defined differently throughout all systems. Participants discussed whether any general fairness rules can be applied to all models. A clear and globalized definition of fairness would clarify the metrics used to measure it, and improve the fairness and user trust of a system.
- **Explainable and Trustworthy Technology:** Trust of a system is difficult to measure, as there are no clearly defined and widely accepted metrics to do so. We discussed how much information users need to understand a system, and how AI should be made explainable. AI should be able to explain its decision making similarly to how a doctor explains a condition to a patient. A doctor does not need to explain down to the cellular level why a patient may be suffering from a condition, but the doctor does need to explain what may have occurred for this condition to be present in the patient (lifestyle choices, hereditary conditions, etc.). Similarly, an AI program does not need to explain why a certain weight was assigned to each node, but rather why an accumulation of weighted nodes is sufficient to make a certain decision. While some researchers have historically argued that people don't require explanations but just want outcomes from AI, sometimes it is required especially to build trust in a system. There is a need for an established formal process to determine the trustworthiness of a system.

Breakout 4

Moderator: George Lan

Notetaker: Jai Moondra

- **Causality methods across disciplines:**

- Statistics - focus is on estimator properties and optimality.
- CS - focus is on scale and data.
- OR - focus is on engineering models and experimental data.
- An important question is how experiments can be controlled: there may be policy and regulatory questions and network effects in industry that prevent random control experiments. For example, similar prices **must** be shown to everyone in cab-sharing. Another question is how to integrate real data with simulation data? In some situations, historical data can be used to design experiments to generate synthetic data. But certain settings like wind turbines pose difficulties to generate synthetic data.
- **Major challenges in XAI / Causality and possible ways in which a multi-disciplinary approach can help:** There are several examples in social sciences where causality helps. One challenge is that policy makers make decisions hoping to influence outcomes, but they sometimes don't accomplish this because they don't understand or use causality. How can they be helped with causal inference? Practical factors like operational costs, etc. matter. Causal inference in decision making can again be used here. For example, housing. Changing leadership and policy makers can focus on areas other than long-term outcomes (they can base their decisions more on politics). Educate them that politics can be done using causality and decision-making. How do we align policy makers' incentives with producing the best outcomes? Also we must bear the burden of explaining our mechanisms to them.
- **Practical Considerations for Causal Inference Methods:**
 - It is sometimes difficult to conduct experiments (ex: gene editing and other medical situations)
 - Regulatory approvals can delay research
 - Confidence needs to be high – repeatability and reproducibility are key
 - Sometimes correlation is the best that can be shown
 - Need a large and representative sample - sometimes there is the risk of some researchers falsifying data due to the small amount of data available
 - **Chains of causality get complicated quickly.** If A implies B which implies C which implies D, conditional effects of A on D are not so easy to understand.
- **Potential/emerging application areas: XAI and Causality**
 - Policy design
 - Medicine
 - Engineering
 - Cybersecurity
 - Social science - interventions

- Security and intelligence
 - Sustainability - protecting parks and preventing poaching (in Africa for instance)
 - Education
-

Panel C: Robustness/Privacy

Moderated by Bistra Dilkina

Slide Presentations

Bo Li - [Figure B-10: Trustworthy Machine Learning: Robustness, Privacy, Generalization, and Their Interconnections](#)

Kush Varshney - [Figure B-11: Problem-Driven Robustness, Privacy, and Fairness](#)

John Abowd - [Figure B-12: Some Lessons from the 2020 U.S. Census Disclosure Avoidance System](#)

The first presentation was given by Bo Li on the topic of Trustworthy Machine Learning. Bo gave several examples of the dangers of Machine Learning, from hacking attempts, such as the Associated Press hack which crashed the stock market in 2013, to the public's concern with using ML programs, such as biometric recognition being used at airports. The goal to improve these Machine Learning programs, as stated by Li, is to close the trustworthiness gap, which can be accomplished in a number of ways.

Firstly, improved Robustness can increase public trust in an ML program simply by making it functional in more situations. Increased robustness not only increases trust in a program by the public, because the program can react correctly to more unexpected inputs, but it also increases the usefulness and accuracy of the program, since the number of situations the program cannot handle is reduced. Generalization, namely, a program's ability to adapt to new and variable data, can also improve trustworthiness. Finally, increased privacy improves trustworthiness, by reducing the risk of a user's data being accessed by a third party. Dr. Li advocated for a holistic approach to improving trustworthy ML, by tackling these three problems simultaneously and in concert.

The second presentation was by Kush Varshney on the topic of Problem-Driven Robustness, Privacy, and Fairness. Varshney started with describing a few problems in the Healthcare industry, starting with the Patient Protection and Affordable Care Act which changed the landscape of the health insurance market in the United States. Insurance companies had to decide which new markets to enter. Markets are defined by geography, age group and other factors. Some considerations were:

- The desire to enroll low-cost (healthy) people to enroll in their plans
- Preventing from accepting or denying enrollment on an individual basis
- Whether or not to offer plans in well-defined markets
- How to use data-driven decision making for determining whether or not to offer plans in new markets

Some of the challenges in this problem included the fact that insurance companies needed cost and demographic data on people who will enroll in new markets whereas they only have this for those enrolled in existing markets. Given that the target domain is unknown, one needs distributionally robust methods. Varshney described a few approaches to solve this problem, including a game theoretic formulation for invariant risk minimization.

Next, Dr. Varshney discussed the topic of privacy concerns in the context of the Health Insurance Portability and Accountability Act which includes the mandate of privacy protection even for health insurance companies' internal planning purposes. He cited that k-Anonymity is a common mathematical interpretation of the privacy condition and one of the approaches is to use k-member clustering (grouping the records so that the smallest group has at least k elements). He briefly described Distribution-preserving k-Anonymity for transfer learning which is an alternative to standard clustering to allow the resulting data to follow the distribution of the original data.

Varshney also discussed the critical need for algorithmic fairness in the health care insurance industry where often health care cost is used as a very poor proxy for health care needs because it could lead to racial discrimination.

Dr. Varshney ended his presentation with some thoughts on Operations Research + Artificial Intelligence.

- He called for a shift in the way AI problems are done which typically use a problem-driven approach. He recommended the incorporation of model-based approaches and called for cross-fertilization with risk management, probability theory, robust optimization, audio signal processing and game theory.
- He commented that there is a lot of similarity between the notion of Explainability and Robustness – they address the same basic problem with different approaches.
- He posed several interesting questions related to OR and AI: Why do we have many discussions on trustworthy ML but not on trustworthy OR? Is it because trustworthiness is not important to OR practitioners or do they do this under a different name (with robustness, sensitivity analysis, etc.)? Is it due to different audiences that consume these results with the OR users being experts within organizations while ML is deployed to non-technical experts? Or is it because there is more media exposure (negative or otherwise) on AI applications and not as much on OR ones?

The final presentation in this session was by John Abowd. Dr. Abowd spoke about the 2020 Census, during which a total of more than 150 billion statistics from 15GB of total data were captured, and the precautions that were taken to secure the data of those who were interviewed. These precautions were necessary to adhere to the increasing privacy protection regulations in the US as across many other parts of the world. After the 2010 Census collected and processed its data, the tabulations were released to the public. Shortly afterwards, it was discovered that the confidential microdata from that census could be accurately reconstructed from the publicly released tabulations. Geographic identifiers were associated with every microdata point, meaning that those who responded to the census could be very accurately identified across the United States, especially those in rural areas.

To prevent this inadvertent disclosure of confidential data from reoccurring during the 2020 Census, the Census team took a number of precautions, such as adhering to a formal privacy protection framework written for the purpose. In particular, a TopDown Disclosure Avoidance System was developed with strict requirements with regard to formal privacy protections. Details are available in the presentation slide deck included with this report.

Some of the key takeaways from John's presentation are:

- Going from suppression to differential privacy is much easier than going from publishing all the microdata to differential privacy.
- 2020 Census data clients had accuracy expectations that modern privacy protection can't support (the 2010 Census basically released all the microdata, although not intentionally).
- It is safe to forecast that AI applications, particularly in industry, are going to face the same conundrum. For instance, advertising executives are not going to like the privacy-protected models (Conventional AI applications are inherently disclosive.)

Wrap-up of Day 1

Ramayya Krishnan

Throughout Day 1, several big themes were discussed pertaining to the topic of Trustworthy AI, including humans in the loop, preference elicitation, robustness, explainability, privacy, etc. Fairness issues were presented in the context of multiple settings such as fairness in healthcare and multi-stakeholder fairness in rideshare scenarios like Uber/Lyft. Several ideas were presented that were worthy of pursuing collaboratively between the OR and the CS communities. Krishnan called upon the participants to write down sketches of these ideas for future collaboration between the two

groups – so that we could discuss them briefly in the final session on Day 2. The intention was to provide a platform to bring researchers together for joint projects.

The participants had robust discussions regarding the boundary between OR and AI: is it conceivable that the differences between the two areas will vanish in the near future or will they persist due to cultural differences? Some of these differences were already discussed in the first workshop and will need to be overcome for fruitful collaborations to emerge. For example, some barriers stem from the culture of conference papers for CS versus journal articles for OR. Is it possible to borrow ideas from each other's preferred mode to enable quicker dissemination of ideas between the two groups? Are there opportunities to co-locate meetings to bring together the AI and OR communities?

Day 2

John Dickerson kicked off Day 2 with a summary of the presentations on Day 1 and asked Prof. Sven Koenig to share some of his recommendations with regard to opportunities for collaboration.

Sven talked about competitions in computer science which have sparked a great deal of interest among researchers, including graduate students exploring new techniques in the field. Some examples of competitions in CS are:

- Robocup logistics competition conducted by the RoboCup Foundation
- NeurIPS 2020 Flatland Competition
- Several competitions advertised by ICAPS 2021 including
 - Learning to run a power network with trust
 - Automatic reinforcement learning for dynamic job shop scheduling problem
 - The dynamic pickup and delivery problem

Several teams participate in such competitions and often the winning teams employ optimization techniques. Sven posed the question: How can OR create challenge problems through their groups? How can both communities collaborate as teams competing in such competitions?

He recommended that we create challenge problems that would spark interest in such collaborations.

Lavanya Marla proposed a sketch of how we could devise challenge problems that would bring together experts in both fields. Some of her suggestions include:

- Classical OR resource allocation problems which have interesting aspects for both groups:
 - Start with classical OR problem

- Resources have preferences/rules that they state - AI models learn rules and preferences
 - OR model makes decisions with added side-constraints
 - Iterate
- Some examples with rich data arise in hospital staffing, airline crew scheduling, semiconductor manufacturing, transportation, etc. The challenges are to get the right data and provide an appropriate infrastructure for solving such problems.

See the concluding section of this report for some challenge problems that were discussed at the end of the workshop.

Panel D: XAI/Causality

Moderated by Yu Ding

Slide Presentations

Zachary Lipton - [Figure B-13: Adapting Predictors under Causally Structured Distribution Shift](#)

Ruoxuan Xiong - [Figure B-14: Design and Analysis of Panel Data Experiments](#)

Yu Ding - [Figure B-15: Causal Inference in Engineering Applications](#)

The theme of this session was on Explainable AI (XAI) and Causality, although the three talks in the session focused primarily on causality and less on XAI. Admittedly, insights gained from understanding causality helps provide explanation to AI methods. The three talks in this session come from three different angles---one talk was given by a computer science (CS) researcher, one by an operations research (OR) researcher, and one on an application of causal inference by an engineering researcher. Altogether the session provides a balanced view on causality research combining both AI/OR perspectives.

The CS talk by Zachary Lipton raises the issue of adapting predictors under causally structured distribution shifts. The speaker discussed the anatomy of a structured shift problem in the context of domain/environments, structure, visibility, manipulation rules, objective, and statistical capabilities, presented examples of structured shift, such as covariate shift, label shift, or missing data shift (source and target data missing at different rates), and stressed the challenges in handling high-dimensional, arbitrarily non-linear data.

The OR talk by Ruoxuan Xiong discussed the design and analysis of panel data experiments in which conventional A/B testing suffers from network interference or

contamination effect. Such problems commonly exist in the experiments run by certain ride sharing companies for testing whether a new feature would improve driver participation rate. A switchback experiment may help but arbitrary switches between control and treatment run into impracticality constraints; for instance, it is not practical for drivers to see different versions of their app every day. The speaker contemplated on a solution using a linear mixed effect model with integer switching variables that are solved optimally through integer programming.

The application talk by Yu Ding presented a success story in which causal inference methods were used for estimating the effects of certain technical upgrades on wind turbines. Both issues of effect estimation and deciding conditional causal relation were discussed. The classical covariate matching method was tailored for estimating turbine upgrade effect in good accuracy, with the caveat of additional needs for bias reduction. The determination of the conditional causal relation is complicated by the presence of autocorrelation in the time-series data, another type of distribution shift different from but complementing the distribution shifts discussed in the CS talk by Zach Lipton. Causal inference methods are shown to make a great positive impact on an important renewable energy application.

Breakouts for Panel D

Breakout 1

Moderator: Xiao Fang

Notetaker: Catherine Gill

- We discussed how AI and OR can contribute to causal inference from observational data and randomized experiments. One significant challenge of causal inference is the interference effect, which refers to the effect of people affecting each other (e.g., people in the treatment group communicating and affecting those in the control group). This challenge can be addressed by panel data experiments, for which OR methods can help decide optimal experiment parameters (e.g., treatment time) and reinforcement learning methods can be employed to improve treatment efficacy.
- We also discussed high-stakes applications of causality and how AI/OR can help. These include quick decisions such as driving decisions made by autonomous vehicles and longer-term decisions such as college admission decisions. We further discussed the implementation of a policy and how to gather data to understand the causal mechanism underneath the policy. For example, a tax policy is implemented at some sites for testing and survey data is then collected from its affected citizens

to understand the causal mechanism underneath the policy. However, we have difficulties in getting truthful survey data because humans are incentivized to lie in surveys. Hence, there might be research opportunities for AI/OR scholars: how to design surveys that solicit truthful data for causal inference. In some other situations, we have to infer features on the basis of surrogate information; for instance we may look at variables which indicate poverty instead of actual poverty statistics. Are there methods to address this issue?

Breakout 2

Moderator: Abhishek Chakraborty

Notetaker: Haley Griffin

- **Machine Learning and Causality Synergies:** Machine learning methods can be very helpful in guiding decision making, but to make truly informed decisions, we have to understand the causality involved. On the other hand, in cases where you have many confounding factors, some of the ML approaches are better able to deal with those factors; so one can take ML approaches and wrap them in causal methods. AI/OR literature should approach the problem of explainable AI similarly to how a doctor treats a patient; identify the problem (the part of a program which is unexplained), prescribe a treatment (an explanation of the causality), and monitor the recovery (keep an eye on the improved program).
- **Empirical Validation Limited by Data:** There are no guarantees in machine learning outcomes. Researchers can validate them empirically, but it is limited by the data set. It is wrong for researchers to use models outside of the domain they were intended for and extrapolate data that may not be correct. It can work out, but oftentimes it doesn't serve the same purpose. When working with end users and funding organizations, researchers should always check to make sure they can rationalize what they are doing. The researcher might have to point out insights they didn't consider, and by fleshing them out they can figure out if their vision is realistic. The greater the use of observational data the better.
- **Need to be Careful before Attributing Causality:** When you are choosing an independent variable for a causal method you are making a choice, and you have to disregard a lot of variables. Researchers should not claim causation until they have made adjustments throughout their research that prove there is a causal relationship rather than just an association. Eliminating confounding variables in attempting to make a fair world may instead lead to a biased one; adjusting to the world requires careful consideration of causality. This can be key when facing complicated optimization problems. The solutions must be data-driven.

Breakout 3

Moderator: Amanda Coston

Notetaker: Ann Schwartz

- **Causality, XAI, and trust:** We discussed the often complicated relationships between causality, explainable AI, and users' trust in AI. We discussed how explanations can engender trust even when they are not causal, and we also debated how causality may improve the fidelity of explanations. We identified a major risk with XAI: explanations may engender misplaced trust in algorithms, and to guard against this, we discussed the importance of listening to and addressing users' concerns with AI.
- **Grand Challenge:** We identified a possible grand challenge that poses the question: What modalities are needed to empower us to determine whether (or not) to trust a data-driven model? How do we create tools and methods to use for this purpose?

Breakout 4

Moderator: Berk Ustun

Notetaker: Cyrus Hettle

This breakout group primarily discussed various aspects of designing good challenge problems and potential application areas for them, as well as suitable topics and participants for the next workshop.

- **Designing Good Challenge Problems:** We discussed how to design challenge problems that could be used to showcase the value of AI in OR and vice-versa. Participants agreed that designing the kinds of challenge problems would require "more than just identifying a real-world problem." One key issue was to strike a middle ground between the kinds of challenge problems across fields. In OR, challenge problems stem from clients, and may yield insights that are "too specific." In AI, challenge problems are too "stylized," which means they produce insights that are "too abstract." A second issue we discussed was developing measurable and well-motivated evaluation criteria (i.e., metrics that can be used to evaluate solutions along with explanations for why these metrics are useful for a given application).
- **Potential Applications for Challenge Problems:** Potential for good challenge problems include: transportation, lending, and hiring. In particular, transportation problems tend to foster interdisciplinary work, for instance, the large

interdisciplinary teams working on these problems within companies. COVID data, such as the challenges faced when rolling out vaccines, is multi-dimensional and could be a valuable source of problems. In contrast, while the COMPAS data set is widely used in fairness work, it can be politically fraught and challenging to separate from baseline issues with the prison system, and we would prefer something more neutral. COMPAS is a landmark dataset to study algorithmic fairness. This dataset was used to predict recidivism (whether a criminal will reoffend or not) in the US. We discussed an example of work with an Alabama criminal justice organization which was presented as an example to some CEOs with the framing of "although this topic is political, the general problem occurs in lending, hiring, etc."

- **Evaluation Methods for Challenge Problems:** In addition, we discussed evaluation methods for challenge problems, which can be a critical step in design ("90% of challenge formulation work is formulating the metric"). Laying out evaluation criteria very clearly (a strength of OR), including quantifying dimensions and specifying correctness criteria, is beneficial. Even giving logical explanations for the criteria could be as valuable as running a challenge. Doing this for 10-12 different examples could help people relate to various prototype scenarios. We discussed some other frameworks for multi-objective problems, including allowing solvers to choose a subset of the objectives or specifying the objectives, but concealing their relative weights. A comparison was made to kidney transplants, where people propose multiple objectives and debate their relative merits.
- **Fairness Problems.** We discussed two broad categories of fairness problems, of which we saw two very different examples in the earlier sessions. Peter Frazier's work on fairness in ridesharing was a multi-stakeholder problem involving riders, drivers, the community, and regulators. In contrast, Dimitris Bertsimas's scarce resource allocation problem on racial differences in post-acuity care compared two groups with similar objectives, with an interesting tension for how to assign costs for people who don't receive care. The timescale of the impact of decisions can also have an effect (immediate vs. several years) or can stretch out over a period of time (changes in retail price, privacy and fairness issues in the Census). It would be useful to study the commonalities and differences between problems in different settings, as the Privacy Forum has done for similar kinds of harms in different applications.
- **Real-World Impact:** Our session discussed how solving "real-world" problems could have a "real-world" impact. In effect, it can be easy to write papers but hard to have an impact outside of the research ecosystem. We highlighted the need to

engage with a broader community of stakeholders – e.g., domain experts, practitioners, and policymakers. For example, sometimes there are disconnects between lawyers and policymakers, both of whom can have important input into constraints and utility for algorithms. This broader engagement is important as it can help us identify relevant problems and solve them correctly.

- **Some Ideas for the Next Workshop:** We highlighted the need to actively recruit individuals from broader communities for future events (e.g., we could ask experts from the federal government to present problems at the next workshop). Bringing in communities and policy beneficiaries in addition to agencies should also be done, though it adds complications. We discussed some examples of crossovers and what we can take from them (Gale-Shapley, Sheldon Jacobson’s work on TSA Pre-Check, FCC actions) and the idea of putting something in place at the regulatory level to provide a mechanism for evaluating the cost of solving problems, not just the utility. We also discussed fostering effective collaboration with these stakeholders outside of AI and OR. Specifically, we need to work effectively when eliciting their preferences and constraints to develop technical solutions, and describing the benefits/limitations of these solutions (i.e., promote the adoption and responsible use).

Closing Session and Challenge Problems

In our concluding discussions, we brainstormed solutions with the intention of uniting AI and OR researchers and pushing the needle towards a more collaborative practice of both subjects. These challenge problems are discussed in further detail below:

We discussed the disparities in fairness between users of ridesharing services, specifically in the city of New York. Transportation Network Companies (TNC’s) such as Uber, Lyft, and taxi services have been observed to provide less reliable service in historically disadvantaged areas, because these areas usually have lower rider demand leading to longer wait times between trips for drivers. To rectify this inequity, financial incentives need to be put in place to encourage more drivers to travel to these areas. Who pays for these incentives, however, is the main point of contention, whether it be the rideshare riders or the TNC’s themselves. To calculate who optimally ought to pay at any given time, and how to best allocate drivers to riders in this multi-agent setting and balance multiple objectives, we propose a competition be created with teams vying to develop the most efficient, equitable, and trustworthy algorithm combining AI and OR practices.

We also recommend creating a summer program, convening Artificial Intelligence and Operation Research experts to educate Ph.D. students on multidisciplinary training in AI and OR. This program would take place annually, similar to Machine Learning summer workshops which have taken place in the past. The summer program could focus on a number of topics on which research has been done in both AI and OR (such as decision making under uncertainty, local search, etc.) and invite both an AI and an OR speaker for each topic. This way, the Ph.D. students understand the techniques that have been developed in both AI and OR as well as their commonalities and differences. The Ph.D. students could work in interdisciplinary teams to solve a challenge problem that requires the integration of AI and OR techniques, to allow these students to overcome some of the challenges listed above.

We also put forth the idea of a joint AI and OR conference on decision making for intelligent robots. OR often uses its optimization techniques to support human decision makers (in business settings), while AI often focuses on autonomous decision making by agents. Robots, for example, must plan their motions and tasks, both individually and as a team, which are independently complex optimization problems but also need to be coordinated. To create robots which are as efficient, accurate, and coordinated as possible, a convergence of AI and OR expertise will be required.

Finally, we suggest another challenge problem which may demand a joint AI/OR approach. The setting is scheduling of nursing and physician staff in healthcare facilities. Healthcare delivery settings have seen significant fluctuations in demand patterns from pre-COVID times, to during COVID and post-COVID, resulting in increased variability of patterns for emergency care, out-of-hospital care, and clinic care. The classical OR version of this problem is an open-loop problem, that is, resources (staff) are assigned to rosters to meet (deterministic or stochastic) demand. However, in reality the process is often iterative, that is, staff often have changing priority over days or weeks of the roster, along with varying demand, which may require that the schedule be re-calculated dynamically to make it human-friendly. To better predict healthcare demands while creating optimal and fair schedules for healthcare staff is a challenge that could be addressed by AI and OR approaches. To view the [Challenge Problems document](#) from the workshop, please use the link provided.

By pursuing each, or any of the proposed solutions above, we believe we can work towards closing the existing gaps between the fields of Artificial Intelligence and Operations Research. Increased collaboration between researchers in both of these fields over time will lead to an erosion of the disparate lexicons both groups use when referring to the same practices, and will lend new expertise to each practice, from which both fields can benefit.

Appendix A: Pre-Workshop Materials

Figure A-1: Workshop Participants

First Name	Last Name	Institution
John	Abowd	U.S. Census Bureau
Hamsa	Bastani	University of Pennsylvania
Dimitris	Bertsimas	Massachusetts Institute of Technology
Tracy	Camp	Computing Research Association
Abhishek	Chakraborty	Texas A&M University
Theodora	Chaspari	Texas A&M University
Amanda	Coston	Carnegie Mellon University
Tapas	Das	University of South Florida
Maria	De-Arteaga	University of Texas Austin
John	Dickerson	University of Maryland
Bistra	Dilkina	University of Southern California
Yu	Ding	Texas A&M University
Xiao	Fang	University of Delaware
Peter	Frazier	Cornell University
Cyrus	Hettle	Georgia Institute of Technology
Nikhil	Garg	Cornell Tech
Cat	Gill	Computing Community Consortium
Haley	Griffin	Computing Community Consortium
Swati	Gupta	Georgia Institute of Technology
Andrew	Kahng	University of California San Diego
Subbarao	Kambhampati	Arizona State University

Sven	Koenig	University of Southern California
Ramayya	Krishnan	Carnegie Mellon University
Radhika	Kulkarni	SAS Institute, Inc. (retired)
George	Lan	Georgia Institute of Technology
Bo	Li	University of Illinois at Urbana-Champaign
Jing	Li	Georgia Institute of Technology
Zachary	Lipton	Carnegie Mellon University
Daniel	Lopresti	Lehigh University
Kristian	Lum	Twitter
Lavanya	Marla	University of Illinois at Urbana-Champaign
Jai	Moondra	Georgia Institute of Technology
Murat	Omay	U.S. Department of Transportation
Jon	Owen	General Motors
Mark	Riedl	Georgia Institute of Technology
Harrison	Schramm	Group W
Ann	Schwartz	Computing Research Association
Thiago	Serra	Bucknell University
David	Shmoys	Cornell University
Alice	Smith	Auburn University
Berk	Ustun	University of California San Diego
Pascal	Van Hentenryck	Georgia Institute of Technology
Kush	Varshney	IBM Research
Phebe	Vayanos	University of Southern California
Cathy	Xia	Ohio State University
Ruoxuan	Xiong	Emory University

Jerry	Zhu	University of Madison-Wisconsin
-------	-----	---------------------------------

Figure A-2: Workshop Agenda

August 16, 2022 (Tuesday)

07:30 AM	NETWORKING BREAKFAST Conference B
08:30 AM	Welcome and Introductions Conference A
09:10 AM	TBD: Brief comments from funding agencies about opportunities for AI funding Conference A
09:30 AM	Panel A: Fairness Conference A Dimitris Bertsimas, Massachusetts Institute of Technology Maria De-Arteaga, University of Texas at Austin Nikhil Garg, Cornell Tech David Shmoys, Cornell University
10:45 AM	BREAK Outside Conference A
11:00 AM	Breakout A Conference A, C, D, and E
11:45 AM	Report Back A Conference A
12:00 PM	NETWORKING LUNCH Conference B Brief remarks from Pascal Van Hentenryck (AI4OPT) and Andrew Kahng (TILOS)

<p>01:00 PM</p>	<p>Panel B: Human Alignment/HCXAI/HCI Conference A</p> <p>Hamsa Bastani, University of Pennsylvania</p> <p>Peter Frazier, Cornell University and Uber</p> <p>Kristian Lum, Twitter</p> <p>Mark Riedl, Georgia Institute of Technology</p>
<p>02:15 PM</p>	<p>Breakout B Conference A, C, D, and E</p>
<p>03:00 PM</p>	<p>Report Back B Conference A</p>
<p>03:15 PM</p>	<p>BREAK Outside Conference A</p>
<p>03:30 PM</p>	<p>Panel C: Robustness/Privacy Conference A</p> <p>Bo Li, University of Illinois at Urbana-Champaign</p> <p>Kush Varshney, IBM</p> <p>John Abowd, US Census Bureau</p>
<p>04:45 PM</p>	<p>Breakout C Conference A, C, D, and E</p>
<p>05:30 PM</p>	<p>Report Back C Conference A</p>
<p>07:00 PM</p>	<p>NETWORKING DINNER Lure, 1106 Crescent Ave NE, Atlanta, GA 30309</p>

August 17, 2022 (Wednesday)

07:30 AM	NETWORKING BREAKFAST Day 2 Conference B
08:30 AM	Recap Day 1 Conference A
09:00 AM	Panel D: XAI/Causality Conference A Yu Ding, Texas A&M Zachary Lipton, Carnegie Mellon University Ruoxuan Xiong, Emory University
10:15 AM	BREAK Outside Conference A
10:30 AM	Breakout D Conference A, C, D, and E
11:15 AM	Report Back D Conference A
11:30 AM	Bringing it all Together Conference A
12:30 PM	NETWORKING LUNCH Day 2 Conference B
01:15 PM	Report Writing Conference A
02:15 PM	End of Workshop

Appendix B: Workshop Slides

Figure B-1: Murat Omay - Overview of ITS JPO Programs

OVERVIEW OF ITS JPO PROGRAMS:
 1) DATA ACCESS AND EXCHANGES
 2) AI FOR ITS

BRIEFING AT THE CCC ARTIFICIAL INTELLIGENCE/OPERATIONS RESEARCH WORKSHOP II (AUGUST 16, 2022)

MURAT OMAIY
 Program Manager - Data Access and Exchanges

Intelligent Transportation Systems Joint Program Office
 U.S. Department of Transportation

AGENDA

- ITS JPO Data Access and Exchanges Portfolio
- ITS JPO AI for ITS Program
- Q&A

DATA ACCESS AND EXCHANGES PORTFOLIO

MISSION, VISION, GOALS

VISION	MISSION	GOAL
Facilitate effective generation, acquisition, governance, management, and analysis of ITS data and code across all modes to advance the state of multimodal ITS research and enable deployment of operations.	An equitable, safe, effective, resilient, and multimodal transportation network enabled and advanced by integrated and connected ITS technologies, data, and tools.	<ul style="list-style-type: none"> Transform the DOT culture to achieve an ethical, business, pragmatic, value-driven, objective-driven data landscape by investing in robust data governance and stewardship Influence our partners and grant recipients to understand the value of increased capabilities in data governance and management by investing in resources and capability maturity assessment Position the DOT in a leading role in data and AI strategic data-governance and AI research, innovative, and transformative data culture Drive innovation to improve the transportation system by increasing the diversity of uses and reuse of data, as well as by enabling secondary research Make better-informed decisions for the sustainability of our investments by robust evaluation programs that encompass primary and secondary research, as well as commercial products developed based on DOT-funded data Extract knowledge to build evidence to support policy and future research

PORTFOLIO FRAMEWORK

Data Functionality	Technology/Investment Area		
	Established/Core (Data Access, Exchange, and AI)	Maturing	Emerging
Adjusted	<ul style="list-style-type: none"> ITS Data Access and Exchange Program ITS Data Exchange ITS Data Exchange ITS Data Exchange 	<ul style="list-style-type: none"> ITS Data Exchange ITS Data Exchange ITS Data Exchange ITS Data Exchange 	<ul style="list-style-type: none"> ITS Data Exchange ITS Data Exchange ITS Data Exchange ITS Data Exchange
Foundational	<ul style="list-style-type: none"> ITS Data Exchange ITS Data Exchange ITS Data Exchange ITS Data Exchange 	<ul style="list-style-type: none"> ITS Data Exchange ITS Data Exchange ITS Data Exchange ITS Data Exchange 	<ul style="list-style-type: none"> ITS Data Exchange ITS Data Exchange ITS Data Exchange ITS Data Exchange

CURRENT ACTIVITIES – RESEARCH AND DEVELOPMENT

- Development of Departmental Data and AI Strategies
- Gap Analysis for High-Value Datasets
- Acquisition of High-value Datasets
- Data Valuation Index (DVI)
- Capability Maturity Models (internal and external)
- Data Governance, Stewardship, and Preservation Methods
- Data for Artificial Intelligence (AI)
- Data for Innovation
- Research, Operational, and AI Data Integration Initiative (ROADII)
- Transportation Information Management System (TRIMS)

CURRENT DIALOGUES

- Data Strategy Working Group led by USDOT OST-R/BTS
 - Goal: Using the NOFO Language as a steppingstone, develop a robust data culture within the DOT, its partners, and its grantees
 - Objective: Capture data governance and management objectives in a common language to be included as requirements in pre-award (NOFO) and post-award (contract/grant) stages, and operationalize the implementation of associated strategies
- CDGC - Data Ethics and Equity Working Group led by USDOT CDO
- Departmental Data Governance WG led by USDOT OCIO
- Data Science Task Force led by USDOT FHWA
- Artificial Intelligence Intermodal Working Group led by USDOT OST
- WH Office of Science and Technology Policy ML-AI Subgroup led by NITRD
- Open Knowledge Network (OKN) led by NSF
- Data Dialogue Series starting at the ITS World Congress in September 2022

RESOURCES & CONTACT INFORMATION

- Factsheets**
 - ITS Data Access and Exchanges Program: https://www.its.dot.gov/factsheets/pdf/PR ITSDataAccessExchanges_Factsheet.pdf
 - Operational Data Environment (ODE): https://www.its.dot.gov/factsheets/pdf/ITS_PO_CODE.pdf
 - ITS Research Data and DataHub: https://www.its.dot.gov/factsheets/pdf/ITS_ResearchData.pdf
- Data Access and Exchanges Portfolio Links**
 - ITS DataHub: <https://www.its.dot.gov/data/index.htm>
 - ITS CodeHub: <https://www.its.dot.gov/code/index.htm>
 - Secure Data Commons (SDC): <https://www.transportation.gov/data/secure>
 - OSS4ITS: https://usdot.gov/its_alliance/its_spaces/OSS4ITS/overview
 - V2X Hub: <https://github.com/usdot-fhwa-ops/v2x-hub>
 - CARMA: <https://highways.dot.gov/research/operations/CARMA>
 - Work Zone Data Exchange (WZDX): <https://www.transportation.gov/iv/data/wzdx>

To learn more about the Data Access and Exchanges Program, contact:

Murat Omay
Program Manager
U.S. DOT ITS Joint Program Office
Murat.Omay@dot.gov

AI FOR ITS PROGRAM

PROGRAM VISION, MISSION, AND GOALS

VISION	MISSION	GOALS
The AI for ITS Program's vision is to advance next generation transportation systems and services by leveraging trustworthy, ethical AI and Machine Learning (ML) for safer, more efficient, and accessible movement of people and goods.	The AI for ITS Program identifies, develops, implements, evaluates, and coordinates technology and policy research to advance the contextualization and integration of AI (including ML) into all aspects of the transportation system.	<ul style="list-style-type: none"> Accelerate deployment and evaluation of AI for ITS applications Spur innovation of potentially transformative AI for ITS applications Maximize value, making the best use of USDOT resources Facilitate nationwide adoption of trustworthy, ethical AI-driven ITS

AI FOR ITS PROGRAM ORGANIZING PRINCIPLES

- Leverage AI to address critical **multimodal ITS challenges and needs**.
- Promote **accountability** by establishing processes to manage, operate, and oversee implementation.
- Prioritize **security and privacy** of sensitive data.
- Promote **quality, reliability, and representativeness of data** sources and processing.
- Identify precise, consistent, and reproducible **performance measures** that are consistent with program objectives and measure performance.
- Foster reliability and relevance over time through **continuous performance monitoring and assessment** of sustained and expanded use.
- Share **open data**, while protecting privacy, and facilitate **open-source development**, while preserving intellectual property, to increase richness of available data and code and promote innovation.
- Share **lessons learned and best practices** to facilitate reproducibility of applications and accelerate adoption of AI-driven next-generation ITS nationwide.

UNDERSTANDING THE AI FOR ITS MARKET

SOURCES SOUGHT NOTICE (2021)

MOTIVATION	PURPOSE	KEY QUESTIONS
<ul style="list-style-type: none"> Build upon AI for ITS stakeholder outreach and research activities Gather input from public, private, and academic sectors on AI-enabled solutions for ITS Identify opportunities for DOT investments to accelerate AI for ITS innovations to deployment 	<ul style="list-style-type: none"> Solicit feedback on: <ul style="list-style-type: none"> "Deployment-ready" applications that leverage AI for ITS needs DOT role and investment areas to facilitate next generation ITS leveraging AI Existing capabilities for developing and deploying AI for ITS 	<ul style="list-style-type: none"> Identification of deployment-ready AI-enabled ITS applications Experience with developing AI-enabled ITS applications Awareness of proven AI-enabled applications from other domains for adoption/integration Top three roles for DOT investments in AI

INSIGHTS FROM SOURCES SOUGHT RESPONSES (#1 OF 3)

- Application Maturity:** Most of the applications identified are mature; degree of maturity varies. Understanding the scalability of applications across locations and network types would require additional research.
- Type of AI:** Supervised ML techniques, including deep learning and computer vision. These techniques require a significant amount of labeled data for training.
- Benefits Measurement:** Majority of the respondents did not specifically measure quantitative benefits, but notes perceived improvements in process efficiency, accuracy, and transportation system safety and mobility.
- Cost Structures:** All costs were emphasized as highly dependent on data and computing requirements as well as geographic scale (e.g., the number of equipped intersections).

INSIGHTS FROM SOURCES SOUGHT RESPONSES (#2 OF 3)

- Data:**
 - Public domain data (specifically, imagery data) may suffer from a lack of quality control.
 - Individual devices/sensors are increasingly becoming primary data sources.
 - Real-time data are becoming increasingly available, bringing newer challenges with processing.
- Cybersecurity:** Respondents did not conduct security analysis of their AI applications. A few followed their organization's best practices (e.g., data anonymization, encryption, minimizing access).
- Trustworthy and Ethical AI:**
 - Consistently produce outputs that are reasonable, auditable, and explainable.
 - Operate within a data governance structure that is ethical, test constantly for bias, and adhere to federal, state, and local requirements, policies, standards, regulations, and laws.
 - Monitor ethical standards as these do not remain fixed and transform in response to evolving situations (what was acceptable 10 years ago, may no longer be considered ethical).

INSIGHTS FROM SOURCES SOUGHT RESPONSES (#3 OF 3)

- Top Challenges:**
 - Unseen data and model drift over time -- requiring ongoing maintenance of AI applications with a human-in-the-loop, especially for decisions of greater consequence
 - Other key challenges: lack of standards; lack of labeled and clean data; lack of compute power; lack of workforce and expertise in AI; skeptical attitude towards change and new technologies; potential proliferation of unethical AI systems; lack of funding
- Top 3 Government Roles:**
 - Resolve AI-related policy issues (e.g., data governance and data sharing policies)
 - Develop standards to ensure easy access and sharing of data for execution
 - Conduct prototype testing/demonstrations of AI-enabled ITS applications

CHALLENGES TO AI ADOPTION AND IMPLEMENTATION

- AI for ITS Challenges and Potential Solutions, Insights, and Lessons Learned Report (in development):**
 - Based on SSN responses, review of literature, and other market research, and ongoing coordination with deployers and others.
- Key findings across the 12 identified challenges:**
 - These challenges for AI adoption and successful implementation are not unique to ITS.
 - There may be tradeoffs between addressing different challenges.
 - Addressing these challenges is an ongoing exercise.
 - Maintaining a human-in-the-loop is helpful in identifying and mitigating these challenges.



POTENTIAL PROGRAM NEXT STEPS

POTENTIAL MULTI-PRONGED APPROACH

- Goal 1: Accelerate deployment and evaluation of AI for ITS** where it can make an immediate impact on safety, mobility, equity
 - BY focusing on near term, more mature applications for scalable rollout
- Goal 2: Spur innovation** regarding potentially transformative AI for ITS applications
 - BY sponsoring a multi-stage innovation challenge, advancing from promising concepts to an operational prototype to a deployable capability
- Goal 3: Maximize value**, making the best use of ITS JPO and departmental resources
 - BY coordinating complementary efforts across the ITS JPO and other USDOT programs as well as external activities
- Goal 4: Facilitate nationwide adoption** of trustworthy, ethical AI-driven ITS
 - BY demonstrating value, disseminating insights and lessons learned, and facilitating peer exchanges and partnerships to support technology development and knowledge transfer

RESOURCES & CONTACT INFORMATION

- Factsheets**
 - AI for ITS Program Overview: https://www.fhwa.dot.gov/research_areas/emerging_tech/ITSJPO_AIforITS_Program.pdf
 - Potential Application of AI in Transportation: https://www.fhwa.dot.gov/factsheets/pdf/AIforITS_Program_Factsheet_RealWorld_AI_PotentialApps_in_Transportation_D4.pdf
 - AI Scenarios in Transportation for Possible Deployment: https://www.fhwa.dot.gov/factsheets/pdf/AIforITS_Program_Factsheet_RealWorld_AI_Scenarios_in_Transportation_D4.pdf
- AI for ITS Program Publications**
 - Summary of Potential Application of AI in Transportation: <https://rosap.nslb.gov/view/dot/50651>
 - AI Scenarios in Transportation for Possible Deployment: <https://rosap.nslb.gov/view/dot/50752>
 - Plan for AI for ITS Program: <https://rosap.nslb.gov/view/dot/53932>

To learn more about the AI for ITS Program, contact:
Robert (Bob) Sheehan, P.E., PTOE
 Acting Chief of Policy, Architecture, and Knowledge Transfer
 U.S. DOT ITS Joint Program Office
Robert.Sheehan@dot.gov

Figure B-2: Dimitris Bertsimas - Improving on Fairness/Bias

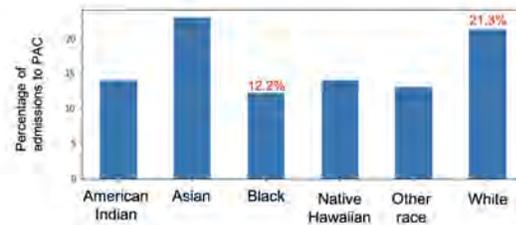
Improving on Fairness/Bias

Hari Bandi and Dimitris Bertsimas

MIT

Based on "The price of Diversity"

Discharge to post acute care



The Problem: Systemic Bias

- Systemic *bias* with respect to *gender, race and ethnicity*, often *unconscious*, but prevalent in datasets involving *choices* made by people.
- Some examples include datasets related to human choices in *college admissions, hiring, lending, or parole* decisions that discriminate against *African-Americans or women*.

The Analytics Edge

Background: Massachusetts General Hospital

- Discharge planning is the development of an *individualized discharge plan* for a patient prior to leaving hospital for home or to a *post acute care (PAC)*.
- Early prediction of PAC needs prior to discharge leads to
 - reducing *hospital length of stay*,
 - unplanned *readmissions*, and
 - improves *patient outcomes*.

Summary

- We propose a novel optimization approach to train classification models on large datasets to *alleviate bias* and *enhance diversity* without significantly compromising on meritocracy.
- Key takeaway: The price of diversity is *low* and sometimes *negative*, that is we can modify our selection processes in a way that *enhances diversity without affecting meritocracy* significantly, and sometimes improving it.

The Analytics Edge

The problem

- The task is to determine discharge disposition for trauma patients within 48-hours after admission.
- Patients are either sent to a *post acute care rehab center* or *home* directly after discharge.
- A successful admission into a PAC depends on
 - Patients' needs,
 - Rehab center agreeing to admit the patient,
 - Patient agreeing to get admitted into a rehab center.

Dataset

- The American College of Surgeons Trauma Quality Improvement Program (ACS-TQIP) database.
- Dataset is sourced from hospitals around the country.
- Features include:
 - patient demographics (age, gender),
 - comorbidities,
 - Emergency Department (ED) vital signs, and
 - injury characteristics (e.g., severity, mechanism).

ML model

- Determine discharge disposition for trauma patients within 48-hours after admission.
- Patients are either sent to a **post acute care rehab center** or **home** directly after discharge.
- Build a Logistic Regression model to predict disposition with AUC = 0.79

Notation

- Each of the patients is assigned an **outcome** $Y \in \{-1, +1\}$ representing either **Entering PAC (+1)** or **not (-1)**.
 - W : set of white patients.
 - B : set of black patients.
 - n_w : total number of white patients.
 - n_b : total number of black patients.
 - p_w : total number of white patients who enter PAC.
 - p_b : total number of black patients who enter PAC.

The Analytics Edge

α -biased dataset

- We call a dataset α -biased if the **difference between the rates of positive observations** among a pair of subgroups W and B based on a protected variable (in this case, race) is at least α .

Definition 1 (α -biased dataset) A dataset $X = \{(x_i, y_i) \mid y_i \in \{-1, 1\}\}$ is said to be α -biased with respect to a pair of subgroups $W, B \subseteq X$ if

$$\left| \frac{\sum_{i \in W} \mathbb{1}(y_i = +1)}{n_w} - \frac{\sum_{i \in B} \mathbb{1}(y_i = +1)}{n_b} \right| \geq \alpha.$$

The Analytics Edge

Demographic parity

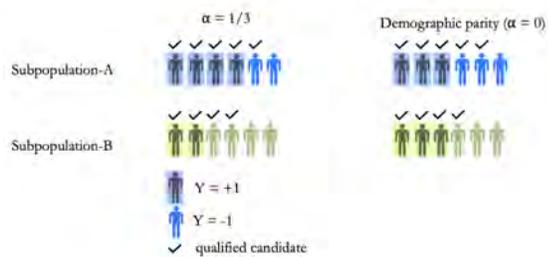
- Demographic parity imposes the condition that a classifier H should **predict a positive outcome** for individuals across groups with **almost equal frequency**.

Definition 2 (Demographic parity) A classifier $H: X \rightarrow \{-1, 1\}$ achieves demographic parity with bias ϵ with respect to groups $W, B \subseteq X$ if and only if

$$\left| \frac{\sum_{i \in W} \mathbb{1}(H(x_i) = +1)}{n_w} - \frac{\sum_{i \in B} \mathbb{1}(H(x_i) = +1)}{n_b} \right| \leq \epsilon.$$

The Analytics Edge

Example of Demographic parity



The Analytics Edge

Proposed solution

- Flip outcome labels (Y) while training your model to achieve demographic parity.
- We propose a Mixed-integer Optimization (MIO) problem that achieves this by introducing binary variables $z_i \in \{0, 1\}$, $i \in [n]$ to decide which outcome labels to flip.

The Analytics Edge

Proposed solution

- The parameters τ_w and τ_b are estimated from the data so that the resulting classifier ensures *demographic parity*.

$$\tau_w \leq \frac{n_b \cdot p_w}{n_w(n_w + n_b)} - \frac{p_b}{n_w + n_b} + \frac{n_b \cdot t}{n_w(n_w + n_b)}$$

$$\tau_b \leq \frac{p_w}{n_w + n_b} - \frac{n_w \cdot p_b}{n_b(n_w + n_b)} - \frac{n_w \cdot t}{n_b(n_w + n_b)}$$

The Analytics Edge

Logistic regression

Logistic regression model

$$\min_{\beta, \gamma} \sum_{i=1}^n \log(1 + e^{-\beta^T x_i + \gamma_i})$$

Parameters of a Logistic regression model

MIO model

$$\min_{\beta, \gamma, z} \sum_{i=1}^n \log(1 + e^{-\beta^T x_i + \gamma_i + z_i t})$$

Binary variables

Updated outcome label

The Analytics Edge

Proposed solution

- If we decide to flip the outcome label of the i^{th} observation: $y_i \in \{-1, 1\}$, the resulting outcome label would be $\bar{y}_i = y_i(1 - 2z_i)$.
- We define a set of n binary variables (z) that flip at most τ_w proportion of labels in W and τ_b proportion of labels in B given by,

$$Z_{\tau_w, \tau_b} = \left\{ z \in \{0, 1\}^n : \frac{\sum_{i \in W} z_i}{n_w} = \tau_w, \frac{\sum_{i \in B} z_i}{n_b} = \tau_b \right\}$$

The Analytics Edge

Logistic Regression

- The dependent variable (Y) is a Bernoulli random variable
 - Y = +1 – “Entering PAC”
 - Y = -1 – “Not entering PAC”
- We seek to predict the probability of a success outcome of the dependent variable Y as a function of independent variables x_1, x_2, \dots, x_k
- We predict the *likelihood* that Y = +1 as follows:
 - $\Pr(Y = +1) = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k}}$ This is guaranteed to be between 0 and 1

The Analytics Edge

MIO model (linearizing product terms)

$$\min_{\beta, \gamma} \min_{z \in Z_{\tau_w, \tau_b}} f(\beta, \gamma) := \sum_{i=1}^n \log(1 + e^{-\beta^T x_i + \gamma_i + z_i t})$$

label flips

$$\sum_{i \in W} z_i \leq \tau_w \cdot n_w$$

$$\sum_{i \in B} z_i \leq \tau_b \cdot n_b$$

Big-M constraints

$$-z_i M_i \leq \gamma_j - \beta_j \leq z_i M_j, \quad i \in [n], j \in [p]$$

$$-(1 - z_i) M_i \leq \gamma_j - \beta_j \leq (1 - z_i) M_j, \quad i \in [n], j \in [p]$$

Implied constraints

$$\sum_{i \in W} \gamma_i = \tau_w \cdot n_w \cdot \beta_j, \quad j \in [p]$$

$$\sum_{i \in B} \gamma_i = \tau_b \cdot n_b \cdot \beta_j, \quad j \in [p]$$

$$z_i \in \{0, 1\}, i \in [n]$$

The Analytics Edge

Additional constraints

- *Maximize likelihood*
- *Demographic parity*
- *Severity of injuries unchanged*
- *Age and gender distribution unchanged*

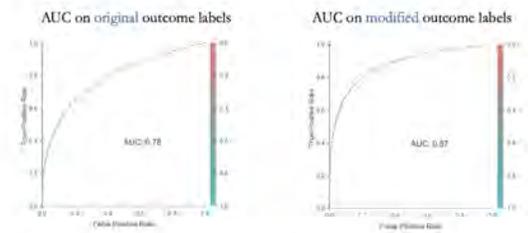
18

Implementation tool

- Train Optimal Classification Trees (OCTs) to provide insights on which attributes of individuals lead to flipping of their labels.
- Construct a dataset based on output of the MIO model. Each defendant is labeled as one of the following:
 - *negative* (patient discharged to home),
 - *high* (patient discharged to PAC), or
 - *no change* (outcome label unchanged)

19

Predictive performance

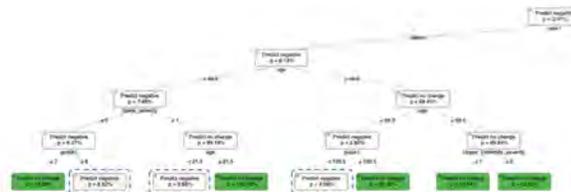


- Alleviating bias improves out-of-sample AUC of OCTs by 8-10%.

21

Implementation tool

Left part of the OCT after splitting on Head severity ≤ 1.0



21

Figure B-3: Maria De-Arteaga - Social Norm Bias: Residual Harms of Fairness-Aware Algorithms

Department of Social Issues
 McCombs School of Business

Social Norm Bias: Residual Harms of Fairness-Aware Algorithms

Maria De-Arteaga, PhD
 Assistant Professor
 Information, Risk and Operations Management Department
 University of Texas at Austin

Algorithmic group fairness

- Measures disparities across **sensitive groups** on a **measure of interest**.

- Examples:
- Gender
 - Race
 - Ethnicity

- Examples:
- Predicted positive
 - False positive rate

2

Algorithmic group fairness

- Measures disparities across **sensitive groups** on a **measure of interest**.

- Examples:
- Gender
 - Race
 - Ethnicity

- Examples:
- Predicted positive
 - False positive rate

- Group fairness constraints enforced via different strategies:
 - pre-processing: transformation of input data
 - in-processing: constraints in optimization objective
 - post-processing: transformation of output scores

3

Algorithmic group fairness

- Very popular due to ease of operationalization, increasingly implemented in ML deployment.
 - Especially true of post-processing (no cost of retraining!)
- Really great that algorithmic fairness research is having an impact! Are we done? Is the problem fixed?



Post-processing

4

Algorithmic group fairness

- Very popular due to ease of operationalization, increasingly implemented in ML deployment.
 - Especially true of post-processing (no cost of retraining!)
- Really great that algorithmic fairness research is having an impact! Are we done? Is the problem fixed?
 - Not really! Risk of reductive definitions and invisibilized harms.



What's Sex Got to Do With Fair Machine Learning?

Algorithmic Fairness from a Non-ideal Perspective

Lay Hoi & Ido Kohler-Hassmann

Stas Fatsopoulos & Zachary C. Lipton

5

What are some of the invisibilized harms?

6

What are some of the invisibilized harms?

Social Norm Bias: Residual Harms of Fairness-Aware Algorithms

Mina Cheng, Maria De-Arteaga, Lester Mackey, Adem Tuzman Kizil
 ICM, Machine Learning for Data Workshop, ICM, Socially Responsible Machine Learning Workshop



7

What are some of the invisibilized harms?

Social Norm Bias: Residual Harms of Fairness-Aware Algorithms

Mina Cheng, Maria De-Arteaga, Lester Mackey, Adem Tuzman Kizil
 ICM, Machine Learning for Data Workshop, ICM, Socially Responsible Machine Learning Workshop



Social Norm Bias (SNoB): The associations between an **algorithm's predictions** and **adherence to social norms**

8

Description and Prescription: How Gender Stereotypes Prevent Women's Ascent Up the Organizational Ladder

Madeline E. Heilman*
Ave 2014 Leadership

How Women Engineers Do and Undo Gender: Consequences for Gender Equality

Madeline E. Heilman*
Ave 2014 Leadership

Gender stereotypes and workplace bias
Madeline E. Heilman

Description and Prescription: How Gender Stereotypes Prevent Women's Ascent Up the Organizational Ladder

Madeline E. Heilman*
Ave 2014 Leadership

How Women Engineers Do and Undo Gender: Consequences for Gender Equality

Madeline E. Heilman*
Ave 2014 Leadership

Gender stereotypes and workplace bias
Madeline E. Heilman

Do algorithms exhibit social norm bias?
Is this bias mitigated by group fairness approaches?
In other words, am I more likely to benefit from group fairness approaches if I sound/act "like a man"?

Experiments using task/dataset *Bias in Bios*

Enter the bio

She is a fifth year PhD student in the joint Machine Learning and Public Policy program at Carnegie Mellon University's Machine Learning Department and Heinz College. She is co-advised by Prof. Arjun Guha and Prof. Alexandra Chouldechova, and she is part of the AuthorLab.

Currently, her main focus is algorithmic fairness, studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support. She is passionate about developing novel machine learning algorithms that are

Abstract Only

SAVE HQ

She is a fifth year PhD student in the joint Machine Learning and Public Policy program at Carnegie Mellon University's Machine Learning Department and Heinz College. She is co-advised by Prof. Arjun Guha and Prof. Alexandra Chouldechova, and she is part of the AuthorLab. Currently, her main focus is algorithmic fairness, studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support. She is passionate about developing novel machine learning algorithms that are motivated by existing policy problems, and understanding how machine learning can better help us overcome important societal challenges. Prior to graduate school she received her B.S. in Computer Science from Universidad Nacional de Colombia and worked as a research intern at Google. She is the recipient of a Fulbright Fellowship.

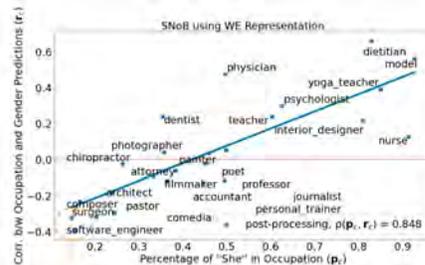
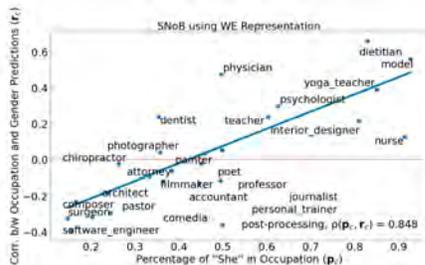
teacher

Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting (FAccT 2019)
M. De Armentia, A. Roemer, H. Wilczk, J. Chirico, C. Borg, A. Chouldechova, S. Ghosh, K. Kimbrough, A. Kulkarni

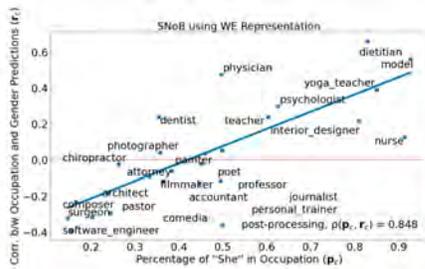
Measuring Inferred Social Norms

- Trained a classifier $G(x)$ to predict group membership (e.g. "he"/"she") from bio x
 - Gender-balanced within occupation
- Use probability of "he"/"she" as measure of aligning to inferred masculine/feminine gender norms
 - Does the biography align with what the algorithm associates as masculine/feminine?
 - Validation to ensure associations learned by algorithms align with human's gender norm associations

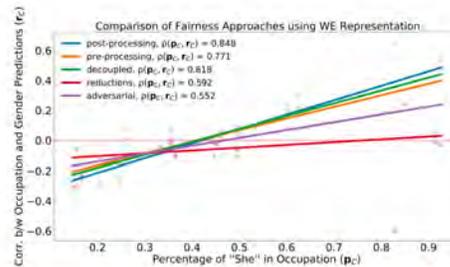
Cryan, Jirina, et al. "Detecting Gender Stereotypes, Lexicon vs. Supervised Learning Methods." ACM CHI, 2020.



women surgeons and software engineers are less likely to be predicted as such if their bio "sound like a woman's"



Strength of correlation increases with occupation's gender imbalance



Gender correlations persist using "fair" interventions

Comparison of Fairness Approaches using WE Representation

post-processing, $p(p_i, r_i) = 0.848$

Invisibilized harms of group fairness

- Algorithms may display **social norm bias**: occupation classification associated to individuals' adherence to social norms.
- Group fairness approaches may **veil harm** to individuals, e.g. fewer opportunities for people perceived as feminine in predominantly masculine occupations.
- Fairness-aware algorithms are not all equal: post-processing that preserves within-group order **does not mitigate** this harm at all.

Comparison of Fairness Approaches using WE Representation

post-processing, $p(p_i, r_i) = 0.848$

Invisibilized harms of group fairness

- Algorithms may display **social norm bias**: occupation classification associated to individuals' adherence to social norms.
- Group fairness approaches may **veil harm** to individuals, e.g. fewer opportunities for people perceived as feminine in predominantly masculine occupations.
- Fairness-aware algorithms are not all equal: post-processing that preserves within-group order **does not mitigate** this harm at all.

Thank you!

Figure B-4: Nikhil Garg - Auditing and Designing for Equity in Government Service Allocation



Auditing and Designing for Equity in Government Service Allocation

Nikhil Garg, Cornell Tech (ngarg@cornell.edu)

Special thanks to Urban Tech Hub at Cornell Tech & NYC Department of Parks and Recreation

Government service allocation

Local government manages many services
 ~8k miles of streets in NYC
 ~700k trees lining streets in NYC
 Housing, sanitation, transportation, etc.



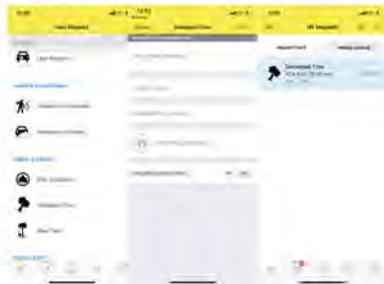
Street trees on Upper East Side in NYC

Operational tasks
 [Learning] What problems are there?
 [Allocation] Which ones to address?
 [Auditing] Did we do a good job?

Desiderata: Efficiency & Equity

311 (crowdsourcing) systems

Cities have a phone number & app to complain to the local government



NYC's 311 system received about 2.7 million requests 2021

These are the primary way the government learns about problems

Pipeline: from incident to work orders



Why is this hard? Uncertainty, heterogeneous + strategic behavior, distribution shifts over time, capacity constraints, pipelined decisions
Research agenda: Audit and improve process's efficiency and equity
Existing collaboration: NYC Department of Parks and Recreation

Understanding reporting behavior



Why? If there are disparities in *who reports* problems, there will be disparities in *what work gets done*

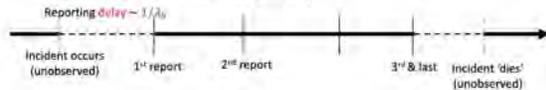
"Equity in Resident Crowdsourcing: Measuring Under-reporting without Ground Truth Data" w/ Zhi Liu (ACM EC 2022)

Model + Method summary

How long does it take for an incident of type θ to be reported?

(Hard because we never observe anything before the first report)

Step 1: Write down a system model where the estimand corresponds to some identifiable quantity



Model + Method summary

How long does it take for an incident of type θ to be reported?

Step 1: Write down a system model where the estimand corresponds to some identifiable quantity

Step 2: Computationally + statistically tractable estimation

$$\# \text{ reports}(i) \sim \text{Poisson}(\lambda_{\theta} \times (b_i - a_i))$$

Spatial smoothing: ICAR Model [Morris et al. 2019]
 Type θ contains an indicator for census tract (2000+ in NYC)
 Then, a_i for each tract is drawn with mean of a_j of neighboring tracts

Results

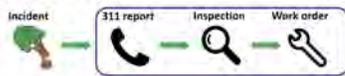
Efficiency: Reporting rates higher for more urgent incidents

Equity: Reporting rates vary substantially by neighborhood

	Manhattan	Queens
High risk hazards	2.5 days	4.7 days
Medium risk tree damage	15 days	28 days
Low risk minor issue	112 days	209 days



Auditing agency decisions in entire pipeline



Questions: Is the agency inspecting the right reports? Are decisions efficient & equitable?

Challenge: Modeling capacity-constrained decisions under uncertainty

Method

- (1) Use ML techniques to estimate incident risk given report characteristics
- (2) Compare "optimal" set allocation decisions with empirical ones

"End-to-end Auditing of Decision Pipelines"
w/ Benjamin Laufer and Emma Pierson

Improving agency decisions



Question: Can we "optimally" re-prioritize inspections and work orders?

Challenge: Want "simple" policies that don't require maintaining an ML model

Method

- (1) Use ML techniques to [robustly] estimate incident risk given report characteristics
- (2) Come up with "service level agreements" for how quickly to address reports

"Making Inspection Decisions: Designing service level agreements"
w/ Zhi Liu

Discussion: ML + Operations



Fairness requires both:
We want to compare decisions for comparable incidents/people/groups



Analyze individual incidents
(characterizing uncertainty, representing data)

Make global decisions
(comparing incidents, allocation under capacity constraints, modeling incentives)

Another example: Recommendation systems

Old school ML view: Predict match between single item and user pair

But there are many global properties of recommender systems

- How users/items affect each other [competition effects]
- How users affect what items are produced [supply-side equilibria]
- How can we recommend sets of items [diverse recommendations]

Joint work with: Christian Borgs, Wenshuo Guo, Meena Jagadeesan, Michael I. Jordan, Karl Krauth, Lydia Liu, Laura Mitchell, Jacob Steinhardt, Gourab K Patro, Lorenzo Porcaro, Qiuyue Zhang, Meike Zehlike

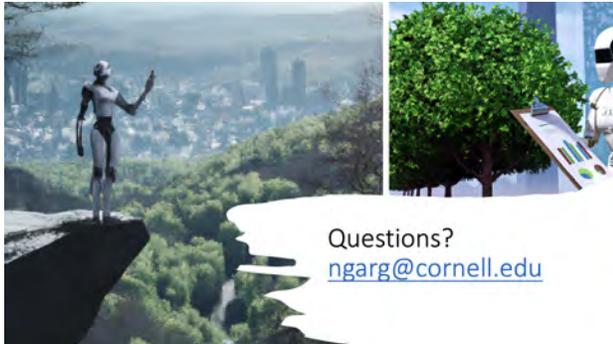
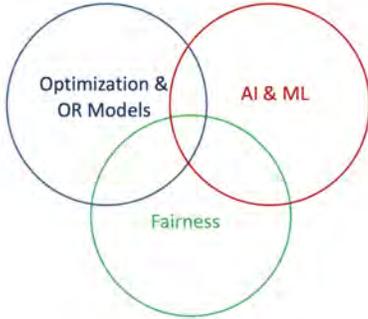


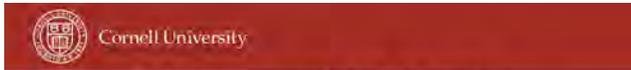
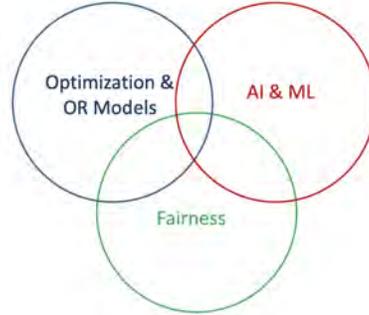
Figure B-5: David Shmoys - Fairness as the Objective in Congressional Districting



Two Examples of this Intersection



Fairness as the Objective

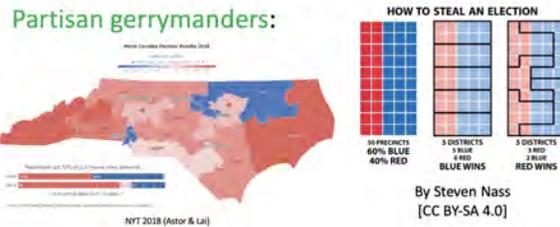


Fairness as the Objective Congressional Districting

- Use historical data to learn distribution for partisan split in any proposed district
- Choose a fairness measure such as “efficiency gap” or “proportionality gap”
- Optimize choice of statewide districting plan to optimize fairness

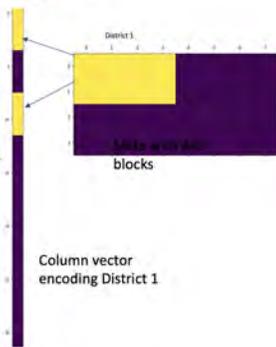


Intentional gerrymandering



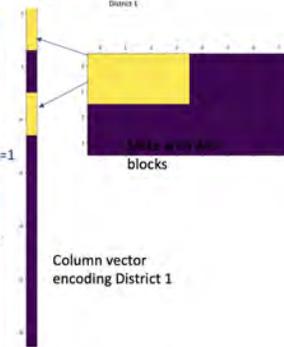
An Idealized Integer Program

- Consider enumerating each possible district
- Each district is subset of blocks in the state
- Encode district as a 0-1 column vector
- In example, yellow = 1, purple = 0
- Let D_1, D_2, \dots, D_N be an enumeration of all *feasible* districts
- Each D_j is a subset of the set of blocks B
- Let $A=(a_{ij})$ be an $|B|$ by N 0-1 matrix where $a_{ij}=1$ if block b_i is in district D_j (& =0 o.w.)
- Introduce 0-1 decision variable x_j where =1 means D_j is one of the districts (& =0 o.w.)
- “Packing constraint”: $\sum_{j=1}^N a_{ij}x_j = 1$, for each $b_i \in B$
- k districts: $\sum_{j=1}^N x_j = k$



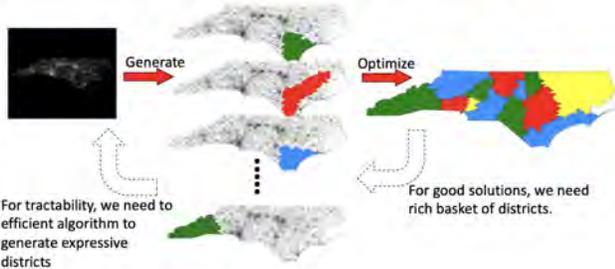
An Idealized Integer Program

- Consider enumerating each possible district
- Each district is subset of blocks in the state
- Encode district as a 0-1 column vector
- In example, yellow = 1, purple = 0
- Let D_1, D_2, \dots, D_N be an enumeration of all *feasible* districts
- Each D_j is a subset of the set of blocks B
- Let $A=(a_{ij})$ be an $|B|$ by N 0-1 matrix where $a_{ij}=1$ if block b_i is in district D_j (& =0 o.w.)
- Introduce 0-1 decision variable x_j where =1 means D_j is one of the districts (& =0 o.w.)
- “Packing constraint”: $\sum_{j=1}^N a_{ij}x_j = 1$, for each $b_i \in B$
- k districts: $\sum_{j=1}^N x_j = k$

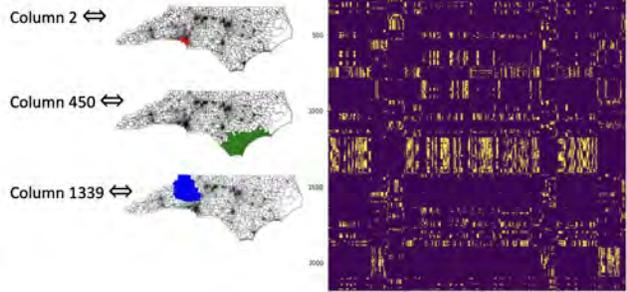




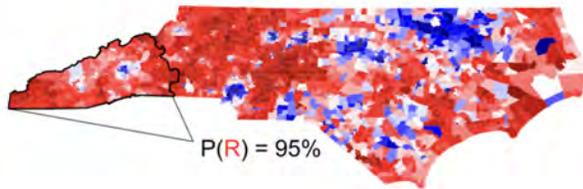
Decoupled Formulation Sketch



Block District Matrix



Political Estimates



Fairness as an Objective

- With traditional formulations, objective function must be linear function of blocks
- With decoupled formulation, objective function must be (piecewise) linear function of districts
 - Coefficient can be arbitrary function of blocks
- Idea: minimize difference between statewide seat share and vote share

Affiliation model: $\nu_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ $\psi_i \sim B(P(\nu_i > .5))$

Linearity of expectation: $E[\sum_{i=1}^k \nu_i - \psi_i] = \sum_{i=1}^k \mu_i - (1 - \Phi(\frac{\mu_i - .5}{\sigma_i}))$

Works for any mapping of votes to seats: $h(\nu_i) = \nu_i$



Optimization – The IP Formulation Again

Input: List of columns (districts)

Output: Indices of the k districts in optimal plan

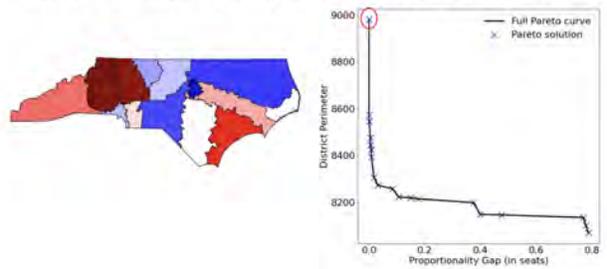
$$c_j = \mathbf{E}(h(\nu_j) - \psi_j)$$

District	Ideal	Seat
j cost	seat	share

$$\begin{aligned} & \text{minimize} && \sum_{j \in D} c_j x_j && (1) \\ & \text{s.t.} && \sum_{j \in D} a_{ij} x_j = 1 && \forall i \in B && (2) \\ & && \sum_{j \in D} c_j x_j = k && (3) \\ & && x_j \in \{0,1\} && \forall j \in D && (4) \end{aligned}$$

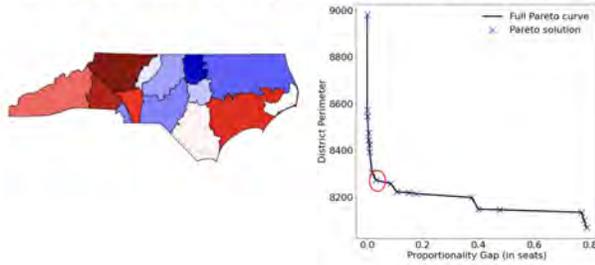


Compactness vs Proportionality Gap Tradeoff

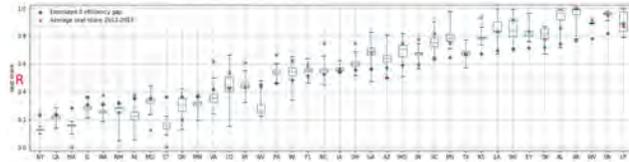




Compactness vs Proportionality Gap Tradeoff



Seat Share Possible Outcomes



- 43% - 62% Expected Democratic Seat Share
- 32-106 Expected seat swaps
- Only 15 states admit proportional plan



This bill requires (1) that ranked choice voting . . . be used for all elections for Members of the House of Representatives, (2) that states entitled to six or more Representatives establish districts such that three to five Representatives are elected from each district, and (3) that states entitled to fewer than six Representatives elect all Representatives on an at-large basis

—Fair Representation Act, H.R. 4000, 2019

We show that 2- or 3-member districts with STV are enough to both inhibit partisan gerrymanders and eliminate natural gerrymanders, without sacrificing “representative” democracy



- [Wes Gurnee, DBS] **Fairmandering: A Column Generation Heuristic for Fairness-Optimized Political Districting.** Proceedings of the 1st SIAM Conference on Applied and Computational Discrete Algorithms (2021).
- [Nikhil Garg, Wes Gurnee, David Rothschild, DBS] **Combating Gerrymandering with Social Choice: the Design of Multi-member Districts.** EC 2022: The 23rd ACM Conference on Economics and Computation, 560-561

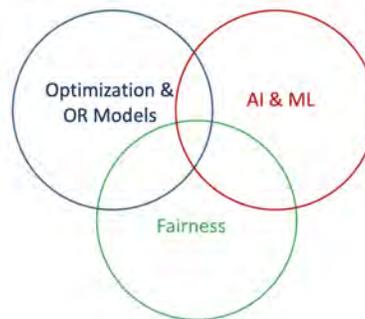


Fairness as the Objective Congressional Districting

- Use historical data to learn distribution for partisan split in any proposed district
- Choose a fairness measure such as “efficiency gap” or “proportionality gap”
- Optimize choice of statewide districting plan to optimize fairness



Fairness as a Constraint





Fairness as a Constraint Managing Ride-Share

- Use historical data to learn rider demand, driver supply, transit times
- Optimize pricing and dispatch policy to maximize long-term profit
- Ensure these policies are “fair”



Fairness as a Constraint Managing Ride-Share

- Use historical data to learn rider demand, driver supply, transit times
- Optimize pricing and dispatch policy to maximize long-term profit
- Ensure these policies are “fair”
- But what is “fair”?



Fairness as a Constraint Managing Ride-Share

- Use historical data to learn rider demand, driver supply, transit times
- Optimize pricing and dispatch policy to maximize long-term profit
- Ensure these policies are “fair”
- But what is “fair”?
- Between robustness & differential privacy, need a notion of optimizing indifference



THANK YOU!

Figure B-6: Hamsa Bastani - Decision-Aware Reinforcement Learning

Decision-Aware Reinforcement Learning

Hamsa Bastani
Wharton School, University of Pennsylvania

Joint work with
Osbert Bastani, Angel Chung, Vahid Rostami, Johnna Sundberg, Musa Komeh, Ashley Schmidt & Lawrence Sandi

Learning and Optimization

- **Goal:** Given response y (e.g., today's demand), compute decision z (e.g., inventory to order) to minimize a known **decision loss** ℓ :

$$z^*(y) = \arg \min_z \ell(z; y)$$
- **Problem:** Optimization parameters y are unknown
- **Strategy:** Predict y based on covariates x (e.g., yesterday's demand)
 - For problems with state, equivalent to model-based offline RL

Learning and Optimization

- **Training phase:** Given examples $\{(x_i, y_i^*)\}$, train a function f_θ to predict y given x :

$$\hat{\theta} = \arg \min_\theta \sum_i \ell(f_\theta(x_i); y_i^*)$$
- **Testing phase:** Given a new x , form prediction $f_\theta(x)$ ("predict") and choose decision $z^*(\hat{y})$ ("optimize")
- **Key question:** What prediction loss $\bar{\ell}(\hat{y}; y^*)$ to use in training?

Learning and Optimization

- **Decision-blind prediction loss:** Use a standard loss such as MSE:

$$\bar{\ell}_{\text{MSE}}(\hat{y}; y^*) = (\hat{y} - y^*)^2$$
- **Decision-aware prediction loss:** Use the decision loss

$$\bar{\ell}(\hat{y}; y^*) = \ell(z^*(\hat{y}); y^*) = \ell(\arg \min_z \ell(z; \hat{y}); y^*)$$
- **Problem:** How to compute $\bar{\ell}$?

Prior Work

- [Elmachtoub & Grigas (2021)] Approximate decision-aware prediction loss $\bar{\ell}$ in the case where:
 - Decision loss ℓ is a linear program
 - Learning restricted to objective (constraint parameters are known)
 - Prediction function $f_\theta(x) = \theta^T x$ is linear
- [Wilder et al (2019a/b)] Heuristics to learn a surrogate loss
- We are interested in:
 - **General decision loss** (nonlinear, unknown constraint parameters)
 - **Any prediction functions** (e.g., f_θ is a random forest)

Approximating the Loss

- We can Taylor expand in \hat{y} around y^* (works well if $\hat{y} \approx y^*$)

$$\ell(z^*(\hat{y}); y^*) \approx \ell(z^*(y^*); y^*) + \nabla_x \ell(z^*(y^*); y^*)^T \nabla_y z^*(y^*) (\hat{y} - y^*)$$
- First term (optimal performance) is constant and can be ignored
- Accounts for:
 - Effect of prediction on decision
 - Effect of decision on decision loss

Training Objective

- Prediction model objective is

$$\arg \min_\theta \sum_i \underbrace{\nabla_x \ell(z^*(y_i^*); y_i^*)^T \nabla_y z^*(y_i^*)}_{\text{constant in } \theta} (f_\theta(x_i) - y_i^*)$$
- Can be interpreted as re-weighting training examples (x_i, y_i^*)
- Compute gradient through **OPT objective** and **OPT decision**
 - Can be computed efficiently (Amos & Kolter, 2017)

Global Health Allocations

1260 health facilities in Sierra Leone

- Severe mismatch & inequality in health supply chains
- 42% of essential medicine needs currently *not* fulfilled

Goal: allocate limited inventory to facilities with highest forecasted need

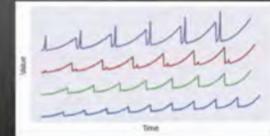
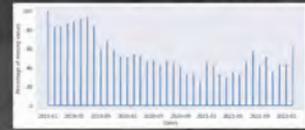
Focal Essential Meds

- **Child Health <5 years of age**
 - Amoxicillin 250mg, Dispersible, Tab
 - Oral Rehydration Salts (ORS), Sachet (correlation to zinc)
 - Zinc Sulphate 20mg, Tab (correlation to ORS)
- **Maternal Health**
 - Oxytocin 10IU, Inj, Amp
 - Magnesium Sulphate 50%, Inj, 10ml, Amp
- **Family Planning (adolescent health, women of child bearing age)**
 - Depot Medroxyprogesterone Acetate (Depo-Provera) 150 mg/ml, Pdr for Inj
 - Ethinylestradiol & Levonorgestrel (Microgynon 30) 30mcg & 150mcg, Tab
 - Jadelle- Levonorgestrel two rod 150mg, Implant

9

Data

- Dhis2, msupply forms
 - Widely used in Mozambique, Cote d'Ivoire, Rwanda, DRC, Chad, etc
- Significant % of missing values → **imputation**
- 1000s of separate time series → **meta-learning**
 - Leverage cross-product, cross-facility correlations
- **Random forest "meta-model"**



10

Out-of-Sample Results



Improve demand forecasts by **34-59%** on held out test set month

11

Stochastic Optimization

- **Decision:** allocations a_n^* across N facilities
- **Objective:** minimize cost of **unmet demand** at each location

$$\xi_n = \max\{\xi_n - s_n - a_n, 0\}$$
 - Current inventory s_n , demand ξ_n
- **Constraints:** fixed budget b , each district cannot hold more than its capacity c_n

$$\xi^{(k)} \geq \xi^{(k)} - s - a$$

$$\xi^{(k)} \geq 0$$

$$s + a \leq c$$
- **Predictions:** draw random demands $\xi_t^{(k)}$ at each facility based on estimated distribution

$$a^* = \arg \min_{a \in \mathbb{R}_{\geq 0}^N} \sum_{k=1}^K \sum_{n=1}^N \xi_n^{(k)}$$

subj. to

$$\sum_{n=1}^N a_n \leq b$$

$$\xi^{(k)} \geq \xi^{(k)} - s - a$$

$$\xi^{(k)} \geq 0$$

$$s + a \leq c$$

* Efficient linear program with sample average approximation

12

For this LP...

- Prediction model objective is approximately

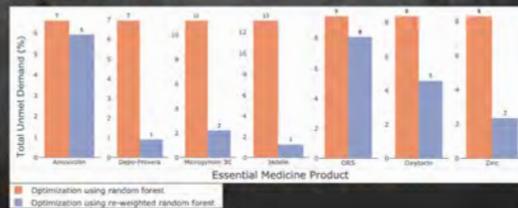
$$\arg \min_{\theta} \sum_{k=1}^K \sum_{n=1}^N \mathbb{I}(\xi_n^{(k)} \geq s_n + a_n) \cdot |f_{\theta}(x_n) - \xi_n^{(k)}|$$

- i.e., we **up-weight** training examples with unmet demand

13

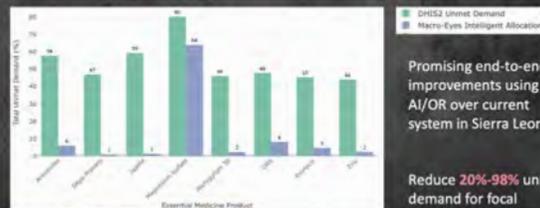
Decision-Aware vs. Decision-Blind

* Compare unmet demand of **decision-blind** vs **decision-aware** random forest + LP for a fixed budget on a held-out test set month



14

End-to-End Relative to Current Approach



Promising end-to-end improvements using AI/OR over current system in Sierra Leone

Reduce **20%-98%** unmet demand for focal essential medicines

Maximum allocation for each product is based on the % of total stock allocated from the Excel tool received for Quarter 1 2022

% of unmet demand = (unmet demand / actual demand) * 100



Thank you!

Questions? hamsab@wharton.upenn.edu



Figure B-7: Peter Frazier - Preference Learning for Stakeholder Management



Preference Learning for Stakeholder Management
Peter Frazier
Cornell & Uber

How testing helped Cornell University become a model of COVID-19 prevention

At the start of the school year, Cornell University implemented a strategy of regular testing and robust contact tracing on campus. The plan was expensive, but it's prevented any major COVID-19 outbreaks at the New York institution.

Cornell COVID-19 Mathematical Modeling Team, Summer 2020:

J. Massey	Ning Duan	Alyf Janmohamed	Brian Liu	Jiayue Wan	Yujia Zhang	Shane Henderson	David Shmoys
-----------	-----------	-----------------	-----------	------------	-------------	-----------------	--------------

Frazier et al., "Modeling for COVID-19 college reopening decisions: Cornell, a case study" PHM 2022

WSJ OPINION

Why Cornell Will Reopen in the Fall

Students will return to their in-classroom. On campus, will get back and better control cases.

... Cornell's decision that it is in the best interest of the health and safety of students, faculty and staff, is to follow a systematic process to directly assess and respond to the needs of the community and to ensure that the university is prepared to respond to any and all potential risks. The school's reopening strategy includes, but is not limited to, regular testing and robust contact tracing on campus. The plan was expensive, but it's prevented any major COVID-19 outbreaks at the New York institution.

... Cornell's decision that it is in the best interest of the health and safety of students, faculty and staff, is to follow a systematic process to directly assess and respond to the needs of the community and to ensure that the university is prepared to respond to any and all potential risks. The school's reopening strategy includes, but is not limited to, regular testing and robust contact tracing on campus. The plan was expensive, but it's prevented any major COVID-19 outbreaks at the New York institution.

... Cornell's decision that it is in the best interest of the health and safety of students, faculty and staff, is to follow a systematic process to directly assess and respond to the needs of the community and to ensure that the university is prepared to respond to any and all potential risks. The school's reopening strategy includes, but is not limited to, regular testing and robust contact tracing on campus. The plan was expensive, but it's prevented any major COVID-19 outbreaks at the New York institution.

Opinion: Morning Editorial Report

We used this graph to help the provost & lab choose testing frequencies

Each letter corresponds to a collection of test frequency assignments

- E.g., "B" means:
 - UG in high-density housing 2x / week
 - Off-campus staff 1x / month
 - Everyone else 1x / week
- Y-axis (Infections) are predicted by a simulation-based from the Cornell COVID-19 mathematical modeling team
- Yellow highlight shows the one chosen

Frazier et al., "Modeling for COVID-19 college reopening decisions: Cornell, a case study" PHM 2022

Successful OR practice requires understanding stakeholders

We took these steps to help the provost & lab choose testing frequencies

- Listen to stakeholders to understand their goals & beliefs
- Run analysis to identify a set of options likely to satisfy them
- Present stakeholders with this set and let them choose

Can AI help with these things?

- Listen to stakeholders to understand their goals & beliefs
- Run analysis to identify a set of options likely to satisfy them
- Present stakeholders with this set and let them choose

Astudillo & F., "Multi-attribute Bayesian Optimization with Interactive Preference Learning", AISTATS 2020
Lin, Astudillo, F., Bakshy, "Preference Exploration for Efficient Bayesian Optimization with Multiple Outcomes" AISTATS 2022

Bayesian Optimization is a black-box derivative-free non-convex optimization method

- x = decision to be made
- Slow computer code or experiment can predict objective function $f(x)$
- $f(x)$ has been evaluated at some x
- Bayesian ML predicts at other x

Acquisition Function: the value of a slow evaluation, e.g., the "expected improvement"

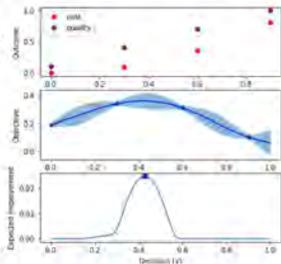
$$EI(x) = E[\max(f(x) - f^*, 0)]$$

where f^* is the best $f(x)$ seen so far.

Bayesian optimization runs a slow experiment at the x with the largest acquisition function value, refits the ML model, & repeats

Sometimes the objective is the utility of an outcome vector

- Slow experiment gives outcome vectors $[h_1(x), \dots, h_n(x)]$
- u = utility function over outcomes
- Objective is $f(x) := u(h(x))$
- Example:
 - $h(x) = [\text{quality}(x), \text{cost}(x)]$
 - $u(h(x)) = \text{quality}(x) - \text{cost}(x)$



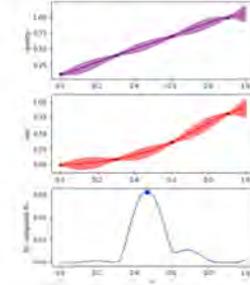
Sometimes the objective is the utility of an outcome vector

You can also put a Bayesian ML model on $h(x)$ and use $E(x) = E[\max\{u(h(x)) - f^*, 0\}]$

Example:

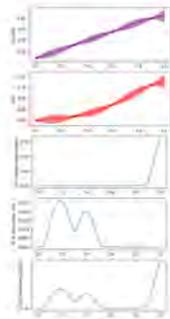
- $h(x) = [\text{quality}(x), \text{cost}(x)]$
- $u(h(x)) = \text{quality}(x) - \text{cost}(x)$

Asatullo & F. "Bayesian optimization of composite functions" ICML 2019...



If the utility function is unknown, we learn it with Bayesian preference learning

- Suppose we don't know the utility function
- Preference learning (for learning the utility fn):
 - Present the stakeholder pairs of outcome vectors.
 - Ask, "Which do you prefer?"
 - Learn a Bayesian ML model on the utility function
 - Active learning decides which pairs to present
- Bayesian optimization (for choosing x to eval):
 - Prediction: Fit a GP to outcomes $h(x)$
 - Acquisition function: $E(x) = E[\max\{u(h(x)) - f^*, 0\}]$ where the expectation E is over the random utility function $u(\cdot)$ and the outcome vector $h(x)$
 - Run the slow computer code where $E(x)$ is largest
- Iterate between asking the stakeholder questions & running the slow computer code

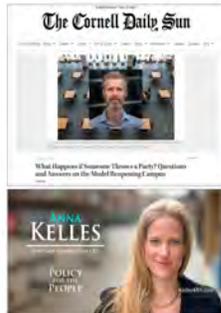


This approach provides lots of value compared to multi-objective optimization

- In use at Meta for product improvement:
 - Stakeholder = product manager
 - Slow experiment = A/B experiment
 - Utility = Quality of Instagram, Facebook, etc.
- Can handle many outcomes; multi-objective optimization usually limited to 3 outcomes
- Can leverage preferences learned from active queries
- Can leverage passive stakeholder observations, e.g., choices in related decision problems



There are many kinds of stakeholder engagement



Lots of room for innovation in AI-enabled stakeholder engagement for OR applications

- Understand stakeholder goals, beliefs & incentives
- Understand how groups of stakeholders influence each other
- Predict how stakeholders will react to communication
- Managing trust (in the OR analyst and her models)

Figure B-8: Kristian Lum - De-biasing "Bias" Measurement

De-biasing "Bias" Measurement

Kristian Lum, Yunfeng Zhang, Amanda Bower

META (Machine Learning Ethics, Transparency, and Accountability), Twitter Inc.

Github: <https://github.com/twitter-research/double-corrected-variance-estimator>



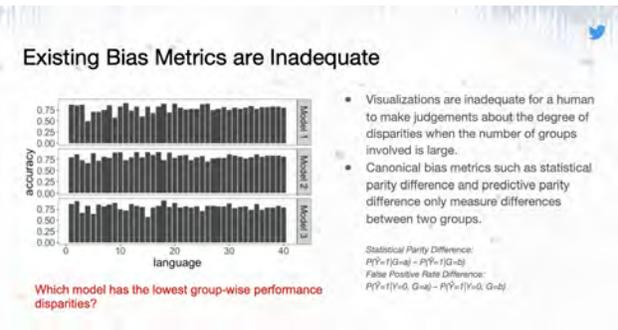
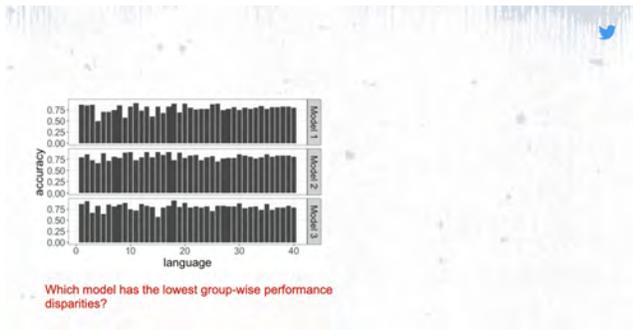
Kristian Lum
Pronouns (She/Her)
@KLDivergence



Yunfeng Zhang
Pronouns (He/Him)
@zywind

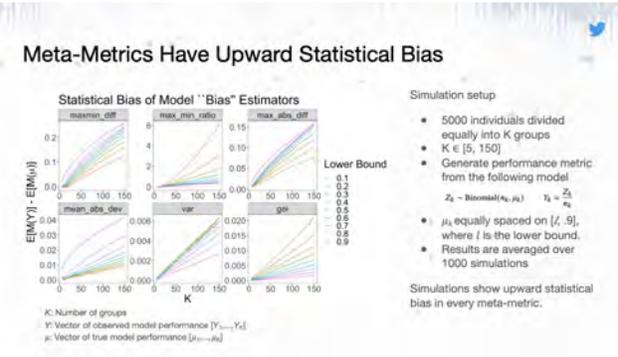


Amanda Bower
Pronouns (She/Her)
@amandarg_



Meta-Metrics

Meta-Metric Name	Formula	Type	Used by
max-min difference	$M_{\text{min-diff}}(Y) = \max_k Y_k - \min_k Y_k$	Extremum	[10]
max-min ratio	$M_{\text{min-ratio}}(Y) = \frac{\max_k Y_k}{\min_k Y_k}$	Extremum	[10, 22]
max absolute difference	$M_{\text{max-abs-diff}}(Y) = \max_k Y_k - \bar{Y} $	Extremum	[2, 32]
mean absolute deviation	$M_{\text{mad}}(Y) = \frac{1}{K} \sum_k Y_k - \bar{Y} $	Variability	[31]
variance	$M_{\text{var}}(Y) = \frac{1}{K-1} \sum_k (Y_k - \bar{Y})^2$	Variability	
generalized entropy index ($\alpha \neq 0, 1$)	$\frac{1}{K(\alpha-1)} \sum_k \left[\left(\frac{Y_k}{\bar{Y}} \right)^\alpha - 1 \right]$	Variability	[9, 44, 46]

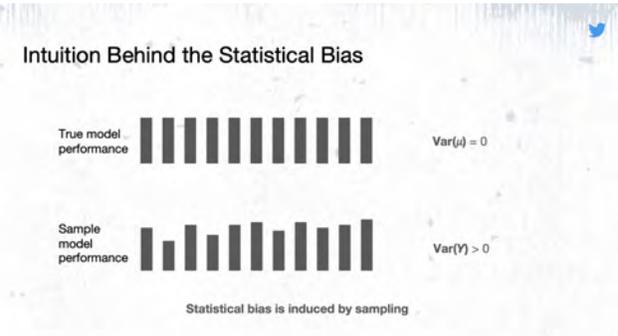


Meta-Metrics Have Upward Statistical Bias

$$\begin{aligned}
 \mathbb{E}[M_{\text{mad}}(Y)] &= \mathbb{E}\left[\frac{1}{K} \sum_{k=1}^K |Y_k - \bar{Y}|\right] \\
 &= \frac{1}{K} \sum_{k=1}^K \mathbb{E}|Y_k - \bar{Y}| \\
 \text{Jensen's inequality} &> \frac{1}{K} \sum_{k=1}^K \mathbb{E}[|Y_k - \bar{Y}|] \\
 &= \frac{1}{K} \sum_{k=1}^K |\mu_k - \bar{\mu}| \\
 &= M_{\text{mad}}(\mu),
 \end{aligned}$$

K: Number of groups
 Y_k : Observed model performance for group k
 μ_k : True model performance $\mathbb{E}[Y_k]$

We can also show mathematically that $M(Y)$ is an upward statistically biased estimate of $M(\mu)$ for several meta-metrics: mean absolute deviation, max absolute deviation, and variance.



Statistically Biased Meta-Metrics Are Problematic

1. The overestimation of the true group-wise performance disparity may cause unnecessary model adjustments. Because of fairness-accuracy trade-off, such adjustments could lead to less performant models for all individuals.
2. Statistically Biased meta-metrics cannot be compared fairly across different grouping methods. Different grouping could lead to different amount of bias.

Correcting for Statistical Bias in The Variance Meta-Metric

$$\hat{M}_{\text{var}}(\mu) = M_{\text{var}}(Y) - \frac{1}{K} \sum_k \hat{\sigma}_k^2$$

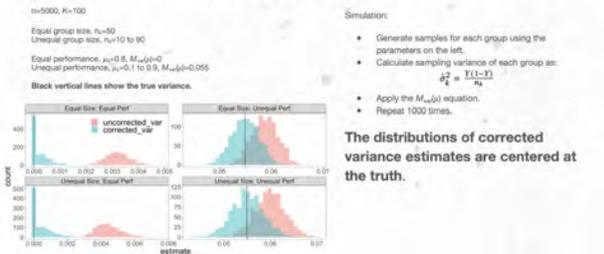
Truncated:

$$\hat{M}_{\text{var}}(\mu) = \max(0, M_{\text{var}}(Y) - \frac{1}{K} \sum_k \hat{\sigma}_k^2)$$

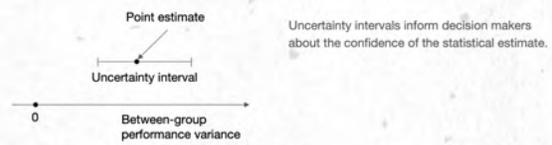
σ_k : Standard error of the observed model performance
 M_{var} : The between-group variance meta-metric.

- Introduced by Cochran (1954) and popularized by Hedge and Olkin (1985) for meta-analysis of between study variability $\hat{\sigma}_k^2$.
- Larger i (e.g. caused by small sample size) leads to larger statistical bias.
- Untruncated version is statistically unbiased when statistically unbiased estimates of the standard errors are available. We follow the convention to use the truncated version to preclude negative variance estimates.

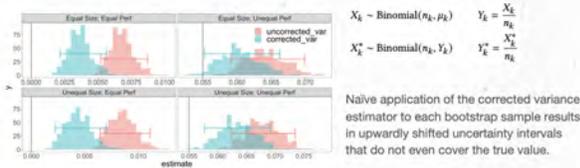
Simulation shows that the correction works



Uncertainty Quantification for MetaMetrics



Uncertainty Quantification By Bootstrapping



Correcting for Statistical Bias in Uncertainty Intervals

Bootstrapping induces variance by itself:

$$\text{Var}(Y_k^*) = \mathbb{E}(\text{Var}(Y_k^* | Y_k)) + \text{Var}(\mathbb{E}(Y_k^* | Y_k))$$

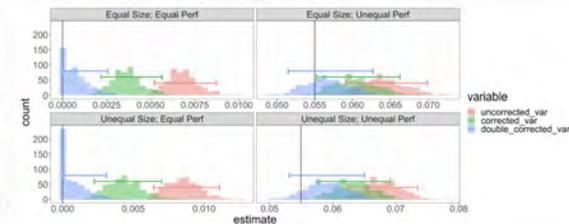
Bootstrap induced var

Sampling var $\hat{\sigma}_k^2 = \frac{Y(1-Y)}{n_k}$

Sampling variance of bootstrapped Y_k^* is:

$$\hat{\sigma}_k^{2*} = \frac{2Y_k^*(1-Y_k^*)}{n_k} - \frac{Y_k^*(1-Y_k^*)}{n_k^2}$$

Corrected Uncertainty Quantification



Corrected Uncertainty Quantification

	uncorrected_var	corrected_var	double_corrected_var
Equal Size; Equal Perf	0.0	0.0	99.7
Unequal Size; Equal Perf	0.0	0.0	99.3
Equal Size; Unequal Perf	15.4	67.6	94.9
Unequal Size; Unequal Perf	10.4	60.4	93.0

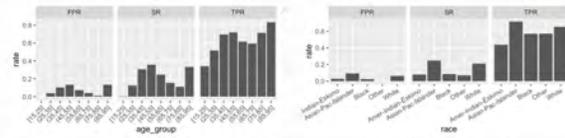
Table 3: Empirical coverage of the 95% bootstrap intervals over 1000 replicates.

Application on The Adult Income Dataset

#	age	workclass	education	education-yr	occupation	relationship	sex
17	Private	Typ	Minor-high school	Mathematics teacher	Married	Male	
18	Private	Self-emp	Minor-high school	Part-time farming	Married	Male	
20	Local-gov	Assoc-acad	Assoc-acad	Production worker	Married	Male	
42	Private	Self-emp	Assoc-acad	Teacher	Married	Male	
19	?	Self-emp	Minor-high school	?	Married	Male	
24	Private	Self	Minor-high school	Other-occupation	Married	Male	
25	?	Self-emp	Minor-high school	?	Married	Male	
43	Self-emp-inc	Prof	Assoc-acad	Prof	Married	Male	
16	Private	Self-emp	Minor-high school	Other-occupation	Married	Male	

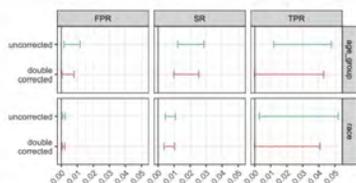
- ~48K individuals from the 1994 census database.
- 14 features
- Label: whether an individual's annual income was above \$50K
- Split the data into 70% train and 30% test
- Trained a gradient-boosted trees classifier.
- 87% accuracy.

Application on The Adult Income Dataset



Age: 8 groups, 15-95 at 10 year interval.
 Race: 5 groups, White, Black, American Indian and Eskimo, Asian and Pacific Islander, and other.
 Performance metrics: False positive rate, selection rate, and true positive rate.

Application on The Adult Income Dataset



The double corrected uncertainty intervals show that in some cases, had we used standard methods, we could have **erroneously** concluded that there are large disparities with statistical confidence.

Contributions and Conclusion

1. We identified meta-metrics for measuring group-wise model performance disparities, particularly in consideration of large numbers of groups.
2. We showed that these meta-metrics are statistically biased measurements.
3. We developed an unbiased estimator for between-group variance based on prior work.
4. We also developed a double-corrected estimator for quantifying the uncertainty of between-group variance.

Future work

- Examine other methods for measuring between-group variability, particularly those from the meta-analysis literature.
- Investigate corrections other than max-min difference.

Caveat: Meta-metrics cannot capture the entirety of the impact of ML systems. Small measured disparities should not be taken as a guarantee that the system is fair or free from adverse impacts.

Figure B-9: Mark Riedl - Toward Human-Centered Explainable Artificial Intelligence

Toward Human-Centered Explainable Artificial Intelligence

Mark Riedl
Georgia Institute of Technology
riedl@cc.gatech.edu
@mark_riedl

Explainable AI (XAI)

- Post-hoc, human-understandable justifications for AI system behavior
- "Opening the black box"

Explainability for Consumer-Facing AI

- Increasingly, those without technical backgrounds in computer science or artificial intelligence are interacting with intelligent systems
- What is *actionable*, and how do we empower non-technical users to achieve remedy?

Human Factors of Explanations

- What effects do different explanation designs have on users?
- Quantitative and qualitative studies of non-AI-expert users and AI-expert users

TL;DR:

- Non-expert values:
 - Contextual accuracy
 - Awareness
 - Strategic detail
 - Intelligibility
 - Relatability
- Unwarranted faith in numbers
- Unanticipated explanatory effects
- Dark patterns & explanation pitfalls
- Do explanations expand the range of actions that consumers have to receive recourse or remedy?

From Trustworthy AI to Trust Calibration

- **System trustworthiness:** aggregate assessment (possibly through verification) of the frequency that a system produces the correct or anticipated behavior
- **Instance trustworthiness:** systems have higher rates of correct/anticipated behavior for some inputs than others
- Question-answering with large neural language models
- Calibrate user's mental model of when the system is right or wrong?

Probing

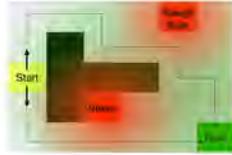
- For multi-hop trivia-like questions, participants can recognize a correct answer 58% of the time
- Reveal internal competencies
- Possibly verifiable by participant knowledge (but at least easier to manually verify)
- Increases participant ability to guess whether the agent's answer is correct to 63% ($p < 0.01$)

Explainable Reinforcement Learning

- Counterfactuals are a consequence of user mental models that do not adequately capture the decision-making of an agent
- Policy doesn't capture why $Q(s, a)$ are the way they are, just that seeking higher values result in larger future expected reward

Experiential Explanations

- Learn to map the influences that certain states have on the utilities of other states
- Present influences as explanations for **why** a particular trajectory is not preferred by the agent's policy
- "I didn't go down because I might fall down the stairs"
- **Actionability:** end-user can update mental model and/or alter the environment to achieve expected behavior



Alabdulkarim and Riedl, Forthcoming

Value Alignment

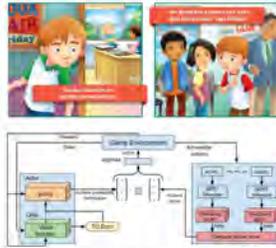
- **Value alignment:** An agent is constrained from performing behaviors that are contrary to human values (often Western values; increasingly consequentialist)
- **Normative alignment:** AI should be biased toward outputs that conform to expected societal and cultural contracts
- How does a system learn "values" or norms? From what data?
- How does a system constrain its behavior according to learned values?

Georgia Tech

Georgia Tech

Normative Alignment

- Learn a normative prior model from stories and educational comics
- Fine-tune language models to avoid generating text about non-normative contexts
- Reinforcement learning to prefer generating normative behavior when trying to achieve a task



Frazer et al. arXiv:1912.03553 Peng et al. arXiv:2001.08764 Al Nahian et al. arXiv:2104.09469

Toward Human-Centered Explainable Artificial Intelligence



Mark Riedl
Georgia Institute of Technology
riedl@cc.gatech.edu
@mark_riedl

Georgia Tech

Georgia Tech

Figure B-10: Bo Li - Trustworthy Machine Learning: Robustness, Privacy, Generalization, and their Interconnections



Trustworthy Machine Learning: Robustness, Privacy, Generalization, and Their Interconnections

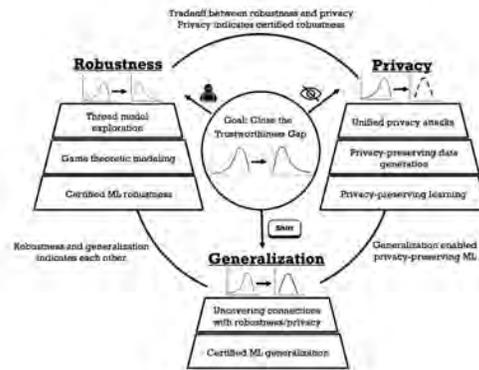
Bo Li

University of Illinois at Urbana-Champaign

Machine Learning in Physical World



Security & Privacy Problems



Lab Test Summary (Stationary)

Adversarial Target:
Stop Sign -> Speed Limit 45
Right Turn -> Stop Sign

Art Perturbation



Physical Adversarial Stop Sign in the Science Museum of London



Possible Vulnerability Disclosure

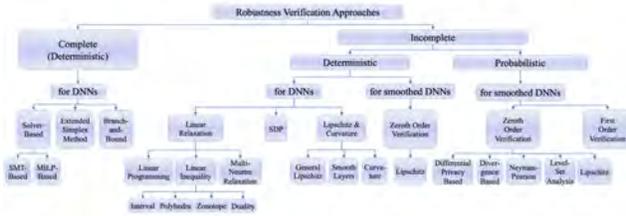
- As of 4/8/21, informed **32 companies** developing/testing AVs
- 12** has replied so far and **have started investigation**



Numerous Defenses Proposed



Certified Robustness For ML



<https://sokcertifiedrobustness.github.io/>

Robust ML Pipeline with Exogenous Information

- Vulnerabilities of statistical ML models: pure *data-driven* without considering exogenous information that cannot be modeled by data
 - Intrinsic information (e.g., spatial consistency)
 - Extrinsic information (e.g. domain knowledge)

COMMUNICATIONS
ACM

Polanyi's Revenge and AI's New Romance with Tacit Knowledge

Research Report, 2021

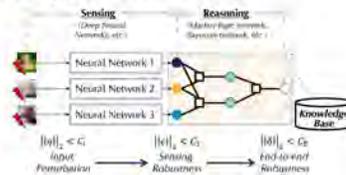


NetHack

The NetHack Challenge: Dungeons, Dragons, and Tourists
eSUNNER 2021 Research

Certified Robustness for Sensing-Reasoning ML Pipelines

- Can we reason about the robustness of an end-to-end ML pipeline beyond a single ML model or ensemble?



- Intuition:** It is hard to attack every sensor in and still preserve their logical relationship
- Goal:** Upper bound the end-to-end maximal change of the marginal probability of prediction
- Challenges:** Solve the minmax for the pipeline

Hardness

Definition 2 (COUNTING). Given input polynomial-time computable weight function $w(\cdot)$ and query function $Q(\cdot)$, parameters α , a real number $\epsilon > 0$, a COUNTING oracle outputs a real number Z such that

$$|Z - \mathbb{E}_{\mathcal{P} \sim \pi} [Q(\mathcal{P})]| \leq \epsilon$$

Definition 3 (ROBUSTNESS). Given input polynomial-time computable weight function $w(\cdot)$ and query function $Q(\cdot)$, parameters α , two real numbers $\epsilon > 0$ and $\delta > 0$, a ROBUSTNESS oracle decides, for any $\alpha' \leq \alpha$ such that $\|\alpha - \alpha'\|_{\infty} \leq \epsilon$, whether the following is true:

$$|\mathbb{E}_{\mathcal{P} \sim \pi} [Q(\mathcal{P})] - \mathbb{E}_{\mathcal{P}' \sim \pi'} [Q(\mathcal{P}')]| < \delta$$

Theorem 4 (COUNTING \leq_t ROBUSTNESS). Given polynomial-time computable weight function $w(\cdot)$ and query function $Q(\cdot)$, parameters α and real number $\epsilon_2 > 0$, the instance of COUNTING, $(w, Q, \alpha, \epsilon_2)$ can be determined by up to $O(1/\epsilon_2^2)$ queries of the ROBUSTNESS oracle with input perturbation $\epsilon = O(\epsilon_2)$.

Theorem 5 (MLN Hardness). Given an MLN whose grounded factor graph is $\mathcal{G} = (V, \mathcal{F})$ in which the weights for interface factors are $w_{ij} = \log(p_i(X_i) / (1 - p_i(X_i)))$ and constant thresholds δ, C , deciding whether

$$\exists \{v_i\}_{i \in V} \text{ s.t. } \forall i, |v_i| < C \text{ and } |\mathbb{E}_{\mathcal{P} \sim \pi} [f_{\text{MLN}}(\{v_i\}_{i \in V})] - \mathbb{E}_{\mathcal{P}' \sim \pi'} [f_{\text{MLN}}(\{v_i\}_{i \in V})]| < \delta$$

is as hard as estimating $\mathbb{E}_{\mathcal{P} \sim \pi} [f_{\text{MLN}}(\{v_i\}_{i \in V})]$ up to ϵ_2 multiplicative error with $\epsilon_2 = O(\delta/C)$.

Robustness of the Reasoning Component

Can we efficiently reason about the provable robustness for the reasoning component when given an oracle for the statistical inference?

Lemma 6 (MLN Robustness). Given access to partition functions $Z_1(\{v_i\}_{i \in V})$ and $Z_2(\{v_i\}_{i \in V})$, and a maximum perturbation C , $\forall \lambda_1, \dots, \lambda_n$ if $\forall i, |v_i| < C$, we have that $\forall \lambda_1, \dots, \lambda_n \in \mathbb{R}$

$$\begin{aligned} & \max_{\{v_i\}_{i \in V}} \ln \mathbb{E}_{\mathcal{P} \sim \pi} [Z_1(\{v_i\}_{i \in V}) + \epsilon_1] \\ & \geq \max_{\{v_i\}_{i \in V}} \left[\sum_{i \in V} \lambda_i v_i \right] + \min_{\{v_i\}_{i \in V}} \left[\sum_{i \in V} \lambda_i v_i \right] \\ & \geq \min_{\{v_i\}_{i \in V}} \ln \mathbb{E}_{\mathcal{P} \sim \pi} [Z_2(\{v_i\}_{i \in V}) + \epsilon_2] \end{aligned}$$

Oracle inference

where

$$\bar{Z}_i(\{v_i\}_{i \in V}) = \ln Z_i(\{v_i\}_{i \in V}) + \epsilon_i + \sum_{j \in V} \lambda_j v_j$$

- When $\lambda_i \geq 0$, $\bar{Z}_i(\{v_i\}_{i \in V})$ monotonically increases w.r.t. v_i ; Thus, the maximal is achieved at $v_i = C$ and the minimal is achieved at $v_i = -C$. When $\lambda_i \leq -1$, $\bar{Z}_i(\{v_i\}_{i \in V})$ monotonically decreases w.r.t. v_i ; Thus, the maximal is achieved at $v_i = -C$ and the minimal is achieved at $v_i = C$.
- When $\lambda_i \in (-1, 0)$, the maximal is achieved at $v_i \in \{-C, C\}$, and the minimal is achieved at $v_i \in \{-C, C\}$ or at the zero gradient of $\bar{Z}_i(\{v_i\}_{i \in V})$ with respect to $v_i = \log \left[\frac{1 - p_i(X_i) \exp(\lambda_i v_i)}{p_i(X_i) \exp(\lambda_i v_i) + 1} \right]$ due to the convexity of $\bar{Z}_i(\{v_i\}_{i \in V})$ in v_i .

Example: PrimateNet (ImageNet)

PrimateNet. The knowledge structure of blue arrows represent the Hierarchical rules between different classes, and red arrows the Exclusive rules. (Some exclusive rules are omitted)



- Hierarchy edge $u \Rightarrow v$:** If one object belongs to class u , it should belong to class v as well.

$$x_u \wedge \neg x_v = \text{False}$$

- Exclusion edge $u \oplus v$:** One object couldn't belong to class u and class v at the same time.

$$x_u \wedge x_v = \text{False}$$

Example: PrimateNet (ImageNet)

Benign Accuracy of models with and without knowledge under different smoothing parameters.

σ	With knowledge	Without knowledge
0.12	0.9670	0.9638
0.25	0.9612	0.9554
0.50	0.9435	0.9371

Certified Robustness and Certified Ratio with different perturbation magnitude C_1 under different smoothing parameters.

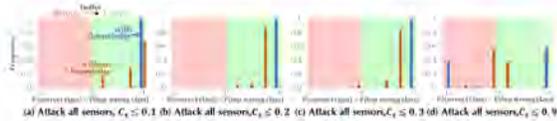
(a) $\sigma = 0.12$					(c) $\sigma = 0.50$				
C_1	α	Cert. Robustness	Cert. Ratio		C_1	α	Cert. Robustness	Cert. Ratio	
0.12	0.05	0.9681	0.9419	0.9529	0.9529	0.05	0.8288	0.7819	0.8287
	0.10	0.8078	0.8095	0.8717	0.9717	0.10	0.7487	0.8488	0.8789
	0.15	0.7308	0.7988	0.8706	0.8706	0.15	0.6907	0.7948	0.8739
0.25	0.05	0.8236	0.8847	0.8706	0.8706	0.05	0.8381	0.8295	0.8633
	0.10	0.7658	0.8428	0.8242	0.8242	0.10	0.7907	0.8446	0.8628
	0.15	0.8220	0.8657	0.8200	0.8300	0.15	0.8285	0.8336	0.8615
0.50	0.05	0.8222	0.8596	0.8200	0.8300	0.05	0.8187	0.8375	0.8612
	0.10	0.8191	0.8824	0.8200	0.8300	0.10	0.8264	0.8318	0.8584

Both the benign accuracy and certified robustness of sensitive reasoning are higher than models w/o knowledge integration.

Example: (NLP) Relation Extraction Task

(NLP) Certified Robustness and Certified Ratio for approaches when all sensing models are attacked.

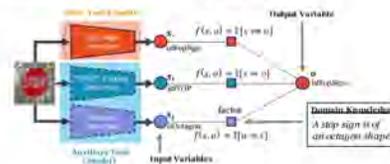
C_S	With knowledge		Without knowledge	
	Cert. Robustness	Cert. Ratio	Cert. Robustness	Cert. Ratio
0.1	1.0000	1.0000	0.9969	0.9969
0.5	1.0000	1.0000	0.9474	0.9474
0.9	0.5882	0.5882	0.3839	0.3839



(a) Attack all sensors, $C_S \leq 0.1$ (b) Attack all sensors, $C_S \leq 0.2$ (c) Attack all sensors, $C_S \leq 0.3$ (d) Attack all sensors, $C_S \leq 0.9$

Example: Knowledge Enhanced ML Pipeline against Diverse Adversarial Attacks

- Example: Robust road sign recognition
- The output of ML models are modeled as input random variables for reasoning
- Permissive knowledge:** s infers γ
- Preventive knowledge:** γ infers s



Knowledge Enhanced ML Pipeline against Diverse Adversarial Attacks

- Lower bound of the pipeline accuracy

Theorem 1 (Convergence of \mathcal{A}^{KEMLP}). For $y \in \mathcal{Y}$ and $\mathcal{D} \in \{\mathcal{D}_b, \mathcal{D}_c\}$, let $\mu_{y, \mathcal{D}}$ be defined as in Lemma 1. Suppose that the modeling assumption holds, and suppose that $\mu_{y, \mathcal{D}} > 0$, for all $K \in \{\mathcal{I}, \mathcal{J}\}$ and $\mathcal{D} \in \{\mathcal{D}_b, \mathcal{D}_c\}$. Then

$$\mathcal{A}^{KEMLP} \geq 1 - \sum_{y \in \mathcal{Y}} \exp(-2\mu_{y, \mathcal{D}}^2/n^2)$$

where v^2 is the variance upper bound to $\mathbb{E}[o] = y|y, W$ with

$$v^2 = 4 \left(\log \frac{1-\alpha_c}{1-\alpha_p} \right)^2 \sum_{K \in \{\mathcal{I}, \mathcal{J}\}} \left(\log \frac{\alpha_{K, \mathcal{D}}(1-\alpha_c)}{\alpha_{K, \mathcal{D}}(1-\alpha_p)} \right)^2$$

$\mu_{y, \mathcal{D}}$ consists of three terms: $\mu_{\mathcal{I}, \mathcal{D}}$, $\mu_{\mathcal{J}, \mathcal{D}}$, and $\mu_{\mathcal{K}, \mathcal{D}}$ measuring the contributions from the main, permissive, and preventative sensors.

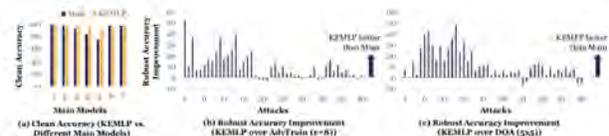
- The accuracy of pipeline is higher than that of the main sensor

Theorem 2 (Sufficient condition for $\mathcal{A}^{KEMLP} > \mathcal{A}^{main}$). Let the number of permissive and preventative models be the same and denote by n such that $n = |\mathcal{I}| = |\mathcal{J}|$. Note that the weighted accuracy of the main model in terms of its truth rate is simply $\alpha_c = \sum_{y \in \mathcal{Y}} \mu_{y, \mathcal{D}} \alpha_{y, \mathcal{D}}$. Moreover, let $K, K' \in \{\mathcal{I}, \mathcal{J}\}$ with $K \neq K'$ and for any $\mathcal{D} \in \{\mathcal{D}_b, \mathcal{D}_c\}$,

$$\mu_{y, \mathcal{D}} = \frac{1}{n-1} \mu_{y, \mathcal{D}} \left\{ \alpha_{y, \mathcal{D}} - 1/2 + \sum_{K \in \mathcal{I}} \mu_{y, \mathcal{D}} - \sum_{K' \in \mathcal{J}} \mu_{y, \mathcal{D}} \right\}$$

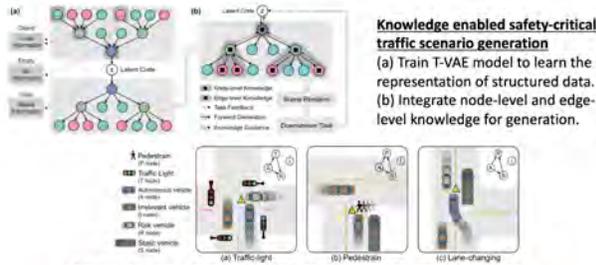
If $\mu_{y, \mathcal{D}} > \sqrt{\frac{1}{2} \log \frac{1-\alpha_c}{1-\alpha_p}}$ for all $\mathcal{D} \in \{\mathcal{D}_b, \mathcal{D}_c\}$, then $\mathcal{A}^{KEMLP} > \mathcal{A}^{main}$.

Knowledge Enabled ML Pipeline Achieves High Benign Accuracy and Robustness under Diverse Attacks



(a) Clean accuracy and (b) (c) robust accuracy improvement of KEMLP ($\alpha = 0.5$) over baselines against different attacks under both whitebox and blackbox settings.

Real-world Case: Autonomous Driving Testing via Logic Reasoning



Causal relationship enabled safety-critical traffic scenario generation
The causal graphs are defined in the upper right for the three scenarios.

The generated safety-critical traffic scenarios can significantly improve the test efficiency of autonomous vehicles.

Testing time certification for RL algorithms



Testing time certification for offline RL algorithms



Unified testing platform for AD via safety-critical scenario generation



Testing time robustness of large language models

AdvGLUE
The Adversarial GLUE Benchmark

What is AdvGLUE?
Advanced GLUE Benchmark (AdvGLUE) is a comprehensive adversarial evaluation benchmark that tests the capabilities of current language models in different adversarial scenarios.

Leaderboard

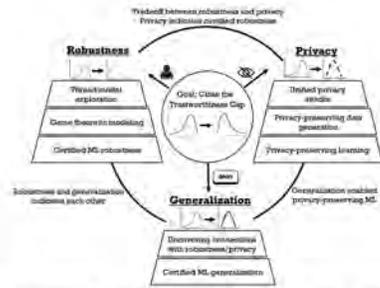
UniFed
A Benchmark for Federated Learning Frameworks

The goal of UniFed is to systematically evaluate the existing federated learning (FL) frameworks. With 11 real-world scenarios, we conduct both qualitative and quantitative performance analysis of existing popular open-sourced FL frameworks. From the perspective of security, confidentiality, communication efficiency, etc., UniFed provides suggestions on framework selection based on the benchmarks. UniFed will point out those frameworks that are more robust. Please check UniFed paper for more details.

Robustness and effectiveness of federated learning frameworks

25

Summary



Closing today's trustworthiness gap requires us to tackle these three grappled problems in a holistic framework, driven by fundamental research focusing on not only each problem but more importantly their interactions.

26

Robust Smart Home

Privacy-Preserving Data Analysis

Topic of Workflow Analysis

Game Theoretic Auditing System for EMR

Large-Scale Auditing Game With Human in the Loop

Robust Learning

Robust Face Recognition Against Poisoning Attack

Privacy Protected Mobile Healthcare

Thank You!
Bo Li
lbo@illinois.edu
<http://boli.illinois.edu/>

Figure B-11: Kush Varshney - Problem-Driven Robustness, Privacy and Fairness

Problem-Driven Robustness, Privacy and Fairness

Kush R. Varshney
Distinguished Research Scientist and Manager
krvarshn@us.ibm.com | @krvarshney

Research IBM

Patient Protection and Affordable Care Act

- Obamacare
- Major provisions came into force in 2014
- Changed the landscape of the health insurance market in the United States
- Health insurance companies had to decide which new markets to enter
 - New markets defined by geography, age group, and other prospect base criteria



IBM Research / August 14, 2012 / © 2012 IBM Corporation

Desiderata of health insurance companies

- Desire low-cost (healthy) people enroll in their plans
- Not allowed to accept or deny enrollment on an individual basis
- Allowed to offer or not offer plans in well-defined markets
- Use data-driven decision making for determining whether or not to offer plans in new markets



IBM Research / August 14, 2012 / © 2012 IBM Corporation

Desiderata of health insurance companies

- Desire cost data on people who will enroll in new markets
- Only have cost data on people who have enrolled in existing markets
- Have demographic data on existing market
- Have demographic data on new market
- Regression problem with covariate shift
 - Also need to consider enrollment: three-population shift



IBM Research / August 14, 2012 / © 2012 IBM Corporation

Regression problem with covariate shift

- Predict cost of individual from demographic features
- Let Y be the cost and X be the features
- Training samples drawn from $p_{X,Y} = p_X p_{Y|X}$
- Test samples drawn from $q_X p_{Y|X}$, where q_X is a different feature distribution
 - Covariate shift
- Minimize risk weighted by importance of each sample under q_X
 - Adaptation
- Take enrollment into account for unique three-population importance weights
 - Fun with Bayes theorem

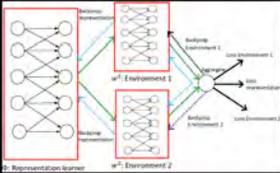
IBM Research / August 14, 2012 / © 2012 IBM Corporation

Adaptation and robustness to distribution shifts

Type	Known Deployment Environment	 covariate shift	 prior probability shift	 concept shift
adaptation	yes	sample weights	obtain labels	
robustness	no	min-max formulation	invariant risk minimization	

IBM Research / August 14, 2012 / © 2012 IBM Corporation

Invariant risk minimization games

IBM Research / August 14, 2012 / © 2012 IBM Corporation

Health Insurance Portability and Accountability Act

- Privacy protection in all releases of health-related information about individuals
 - Even for health insurance companies' internal planning uses
- k-Anonymity is a common mathematical interpretation of the privacy condition
 - Group the records so that the smallest group has at least k elements
 - k-Member clustering
- Standard approach destroys the probability distribution, which is needed for covariate shift adaptation



IBM Research / August 14, 2012 / © 2012 IBM Corporation

Figure B-12: John Abowd - Some Lessons from the 2020 U.S. Census Disclosure Avoidance System

Some Lessons from the 2020 U.S. Census Disclosure Avoidance System

John M. Abowd
 Chief Scientist and Associate Director for Research and Methodology
 U.S. Census Bureau
 Computing Community Consortium, INFORMS, ACM SIGAI
 Artificial Intelligence/Operations Research Workshop II
 Panel C: Robustness/Privacy, Tuesday, August 16, 2022, 3:30pm

The views expressed in this talk are my own and not those of the U.S. Census Bureau. OMS Project ID: P-7502798.
 DRB Clearance numbers: CBDRB-FY20-DSEP-001, CBDRB-FY22-DSEP-003, CBDRB-FY22-DSEP-004.

Bottom Line Up Front:

Going from suppression to differential privacy is much easier than going from publishing all the microdata to differential privacy.

United States
 Census

Forecast:

AI applications, particularly in industry, are going to face the same conundrum. Advertising executives are not going to like the privacy-protected models. (Conventional AI applications are inherently disclosive.)

United States
 Census



United States
 Census

Acknowledgements

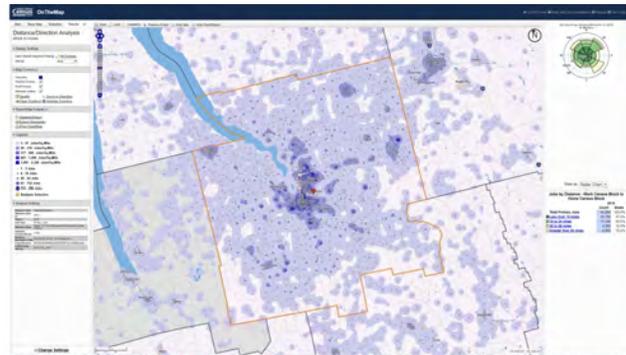
Co-authors on the HDSR paper: Robert Ashmead, Ryan Cumings-Menon, Simon Garfinkel, Christine Heiss, Robert Johns, Daniel Kifer, Philip Leclerc, Ashwin Machanavajjhala, Brett Moran, William Sexton, Matthew Spence, Pavel Zhuravlev
[The 2020 Census Disclosure Avoidance System TopDown Algorithm - Special Issue 2: Differential Privacy for the 2020 U.S. Census \(mit.edu\)](#)
 Annual Review of Statistics paper (with Michael Hawes): [2206.03524] Confidentiality Protection in the 2020 US Census of Population and Housing (arxiv.org)
 This presentation also includes work by the Census Bureau's 2020 Disclosure Avoidance System development team, Census Bureau colleagues, and our collaborators from the following Census Bureau divisions and outside organizations: ADCOM, ADDC, ADRM, CED, CEDA, CEDSC, CES, CSRM, DCMD, DTD, ESMD, GEO, POP, TAB, CDF, Econometrica Inc., Galois, Knexus Research Corp, MITRE, NLT, TI, and Tumult Labs.
 I also acknowledge and greatly appreciate the ongoing feedback we have received from external stakeholder groups that has contributed to the design and improvement of the Disclosure Avoidance System.

United States
 Census

Translation:

2020 Census data clients had accuracy expectations that modern privacy protection can't support (2010 Census basically released all the microdata, although not intentionally).

United States
 Census



Major data products from the 2020 Census:

- Apportion the House of Representatives (April 26, 2021)
- Supply data to all state redistricting offices (August 12, 2021)
- Demographic and housing characteristics (May 2023)
- Detailed demographic and housing characteristics (Part A August 2023; Part B TBD; Supplemental DHC TBD)
- American Indian, Alaska Native, Native Hawaiian data (Included in Part A Detailed DHC; August 2023)

For the 2010 Census, this was *more than 150 billion* statistics from 15GB total data.

United States
 Census

Reconstructing the 2010 Census-I

- The 2010 Census collected information on the age, sex, race, ethnicity, and relationship (to householder) status for 308,745,538 million individuals. (about 1.5 billion confidential data points; Garfinkel et al. 2019)
- The 2010 Census data products released over 150 billion statistics
- Internal Census Bureau research confirms that the confidential 2010 Census microdata can be accurately reconstructed from the publicly released tabulations
- This means that all the tabulation variables for 100% of the person records on the confidential data file can be accurately reproduced from the published tabulations
- Based on Dinur and Nissim (2003) and Dwork and Yekhanin (2008)
- A violation of the 2010 Census contemporaneous disclosure avoidance standards for 2010 Census microdata files



Reconstructing the 2010 Census-II

- A violation of the 2010 Census contemporaneous disclosure avoidance standards for 2010 Census microdata
 - The reconstructed microdata are not a sample; there is one record for every person enumerated in the 2010 Census, and the geographic identifier on that record is always correct (matches the geographic identifier on the confidential input file—the Hundred-percent Detail File, which was swapped)
 - The reconstructed microdata have geography identifiers with an average population of 29 (50, if only occupied blocks are counted)
 - The reconstructed microdata have U.S.-level demographic calls (race, ethnicity) with fewer than 10,000 persons
- The standards for releasing microdata from the 2010 Census required (McKenna, 2019)
 - Sample (10% rate was used)
 - Restrict geographic identifiers to areas with at least 100,000 persons (Public-use Microdata Areas)
 - Collapse demographic categories until the national population in 1-way marginals contains at least 10,000 persons
 - The standards for tabular data permitted universe files, block geography, and low-U.S. population demographic groups (McKenna, 2018) on the assumption that microdata reconstruction was infeasible
- These are the reason the Data Stewardship Executive Policy Committee instructed the 2010 Census not to use swapping as the main protection for 2010 Census products from the reconstruction evidence alone; swapping plus aggregation did not protect the 2010 Census confidential microdata properly



Reconstructing the 2010 Census: What did we find?

Table 1 Agreement Rates (Reconstruction to CEF) by Block Size

Block Size	Total	1-9	10-49	50-99	100-249	250-499	500-999	1,000+
Agreement	91.8%	74.0%	93.0%	93.1%	92.1%	91.3%	90.6%	91.5%

DRB clearance number CBDRB-FY22-DSEP-004; Source: Hawes (2022).

- Block, sex, age (exact/binning in 38 categories), race (OMB 63 categories), and ethnicity were reconstructed:
 - Exactly for 91.8% of the population
 - Exactly 74.0% in the smallest population blocks, but 93.0% in blocks with 10-49 people and 93.1% in blocks with 50-99 people
 - An external user can confirm that these solutions correspond to the exact record in the confidential data for 65% of all blocks using only the published data because there is provably one and only one reconstruction possible in these blocks. That user can identify population uniques on any combination of reconstructed variables.



This is one of the principal failures of the 2010 tabular disclosure avoidance methodology — swapping provided protection for households deemed “at risk,” primarily those in blocks with small populations, whereas for the for the entire 2010 Census 57% of the persons are population uniques on the basis of block, sex, age (in years), race (OMB 63 categories), and ethnicity. Furthermore, 44% are population uniques on block, age and sex. Aggregation provided no additional protection for most blocks.



Table 5 Distribution of Population and Population Uniques by Block Population Size

Block Population Bin	Number of Blocks in Bin	2010 Census Population in Bin	Cumulative Population	Percent of Population in Bin	Cumulative Percent of Population	Population Uniques (block, sex, age) in Bin	Percent of (Block, sex, age) Uniques in Bin
TOTAL	11,078,297	308,745,538				135,432,888	43.87%
0	4,871,270	0	0	0.00%	0.00%		
1-9	1,821,665	8,069,681	8,069,681	2.61%	2.61%	7,670,927	95.06%
10-49	2,671,753	67,597,683	75,667,364	21.89%	24.51%	53,435,603	79.05%
50-99	994,513	69,073,496	144,740,860	22.37%	46.88%	40,561,372	58.72%
100-249	540,455	80,020,916	224,761,776	25.92%	72.80%	27,258,556	34.06%
250-499	126,344	42,911,477	267,673,253	13.90%	86.70%	5,297,867	12.35%
500-999	40,492	27,028,992	294,702,245	8.75%	95.45%	1,051,924	3.89%
1000+	9,805	14,043,293	308,745,538	4.55%	100.00%	156,639	1.12%

DRB clearance number CBDRB-FY21-DSEP-003



Table 2 Reidentification rates for population uniques

Match file	Universe	Putative rate ^a	Confirmed rate ^b	Precision ^c
Commercial	All data defined persons ^d	60.2%	24.8%	41.2%
	Population uniques ^e	23.1%	21.8%	94.6%
CEF	All data defined persons ^d	97.0%	75.5%	77.8%
	Population uniques ^e	93.1%	87.2%	93.6%

^aThe number of records that agree on block, sex, and age (exact/binning), divided by the total number of records in the universe. ^bThe number of records that agree on PIK (the Census Bureau's internal person identifier), block, sex, age (exact/binning), race, and ethnicity, divided by the total number of records in the universe. ^cThe number of confirmed reidentifications (records that agree on PIK, block, sex, age (exact/binning), race, and ethnicity) divided by the number of putative reidentifications (records that agree on block, sex, and age (exact/binning)). ^dAll individuals with a unique PIK identifier within the block (276 million persons for the 2010 Census). ^eAll data defined individuals who are unique in their block on sex and exact/binning age. DRB clearance number CBDRB-FY22-DSEP-004; Data are from Abowd et al. (under review) released in Hawes (2022). Abbreviations: CEF, Census Edited File; DRB, Disclosure Review Board; PIK, Protected Identification Key.



Table 3 Reidentification rates for population uniques of the block's modal and non-modal races

Match file	Population uniques ^a	Putative rate	Confirmed rate	Precision
Commercial	All population uniques	23.1%	21.8%	94.6%
	Of the modal race	25.3%	24.2%	95.3%
	Of the nonmodal races	13.7%	12.2%	89.2%
CEF	All population uniques	93.1%	87.2%	93.6%
	Of the modal race	94.8%	91.3%	96.3%
	Of the nonmodal races	86.2%	70.2%	81.5%

^aIndividuals who are unique in their block on sex and exact/binning age. DRB clearance number CBDRB-FY22-DSEP-004. Data are from Abowd et al. (under review) released in Hawes (2022). Abbreviations: CEF, Census Edited File; DRB, Disclosure Review Board.

- This is not a statistical use, and both the Census Act (13 U.S. Code § 8(b) & 9) and CIPSEA (44 U.S. Code § 3561(11) 'Statistical Purpose') clearly prohibit releasing data that support non-statistical uses.



All 2020 Census Publications

- Will all be processed by a collection of differentially private algorithms (Dwork et al. 2006a, 2006b; Dwork 2006) using the zero-Concentrated DP privacy-loss accounting framework (Bun and Steinke 2016) implemented with the discrete Gaussian mechanism (Canonne et al. 2020, 2021)
- Using a total privacy-loss budget set as policy, not hard-wired, determined by the Data Stewardship Executive Policy Committee
- Production code base, technical documents, and extensive demonstration products based on the 2010 Census confidential data have all been released to the public
- More information: https://www.census.gov/newsroom/blogs/research-matters/2019/10/balancing_privacyan.html



TopDown Algorithm System Requirements

- The 2020 Disclosure Avoidance System's TopDown Algorithm (TDA) implemented formal privacy protections for the P. L. 94-171 Redistricting Data Summary File
- Planned for use in the Demographic Profiles, Demographic and Housing Characteristics (DHC), and Special Tabulations of the 2020 Census
- TDA system requirements include:
 - Input/output specifications
 - Invariants
 - Edit constraints and structural zeros
 - Tunable utility/accuracy for pre-specified tabulations
 - Privacy-loss budget asymptotic consistency
 - Transparency



What is a histogram?

Record ID	Block	Race	...	Sex
1	1001	Black	...	Male
2	1001	Black	...	Male
3	1001	Asian	...	Female
4	1001	Asian	...	Female
5	1001	Black	...	Male
6	1001	AIAN	...	Female
7	1001	AIAN	...	Male
8	1001	Black	...	Female
9	1001	Black	...	Female



Attribute Combination (Block/Race/.../Sex)	# of Records
1001/AIAN/.../Male	1
1001/AIAN/.../Female	1
1001/Asian/.../Male	0
1001/Asian/.../Female	2
1001/Black/.../Male	3
1001/Black/.../Female	2
...	...

Histogram: Record count for each unique combination of attributes (including location), equivalent to the fully saturated contingency table, vectorized, and with structural zeros removed or imposed by constraint

Microdata: One record per respondent



Noisy Measurements

- TDA allocates shares of the total privacy-loss budget by geographic level and by query
- Each query of the confidential data will have noise added to its answer
- The noise is taken from a probability distribution with mean=0, and variance determined by the share of the privacy-loss budget allocated to that query at that geographic level
- These noisy measurements are independent of each other, and can include negative values, hence the need for post-processing



Zero-Concentrated Differential Privacy (zCDP)

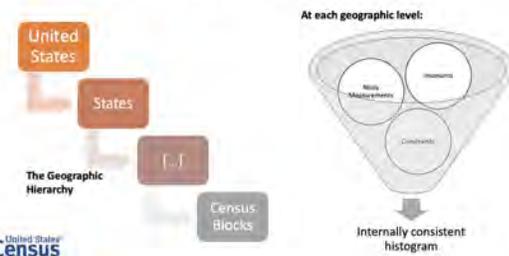
- Privacy-loss parameter: ρ (Bun and Steinke 2016)
- ρ -based privacy-loss budgets can be converted to any single point along a continuum of (ϵ, δ) pairs. Analysis of the privacy protection afforded by a ρ budget should use the entire continuum, not a single (ϵ, δ) point. Some formulas provide tighter bounds on the (ϵ, δ) curve implied by a particular value of ρ . TDA uses this one:

$$\epsilon = \rho + 2\sqrt{-\rho \log_e \delta}$$

- Noise distribution: discrete Gaussian (Canonne et al. 2020, 2021)
- The expected variance of any noisy measurement can be estimated by knowing the total privacy-loss budget and the share of ρ allocated to that query at that geographic level [see Appendix B of [Abowd et al. \(2022\) technical paper](#)]

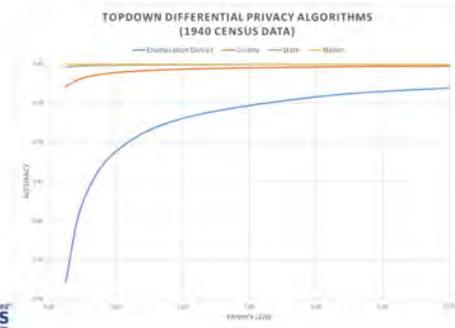
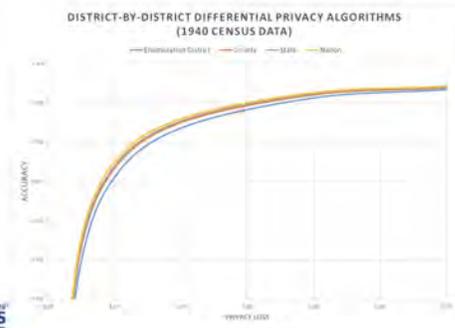


The TopDown Approach



Naïve Method: BottomUp or Block-by-Block

- Apply differential privacy algorithms to the most detailed level of geography
- Build all geographic aggregates from those components as a post-processing



Benefits of TDA Compared to Block-by-block

- TDA is in stark contrast with naïve alternatives (e.g., block-by-block or bottom-up)
- TDA disclosure-limitation error does not increase with number of contained Census blocks in the geographic entity (on spine)
- TDA yields increasing relative accuracy as the population being measured increases (in general), and increased count accuracy compared to block-by-block
- TDA “borrows strength” from upper geographic levels to improve count accuracy at lower geographic levels (e.g., for sparsity)



If you feed TDA 16.6 billion differentially private measurements (23 trillion for DHC), it will do a good job that completely satisfies no one.

(This was predicted in Abowd and Schmutte 2019.)



Accurate, but to whom?

- DAS operates under interpretable formal privacy guarantees, given privacy-loss budgets
- Accuracy properties depend upon the output metric (use case)
- Distinct groups of data users will have a particular analyses they wish to be accurate
- Tuning accuracy for a given analysis can reduce accuracy for other analyses
- Policy makers must consider reasonable overall accuracy metrics for privacy tradeoffs



Deep Dive: Redistricting Data

- Legislative districts for politically defined entities of arbitrary size
- Must be (approximately) equal populations in each district
- Districts must be consistent with Section 2 scrutiny under the 1965 Voting Rights Act
 - Large minority populations cannot be clustered into a few districts
 - Majority-minority districts (approximately 50%+ minority population) must be drawn when feasible
- Focus statistics: total population, ratio largest race/ethnic population to total population

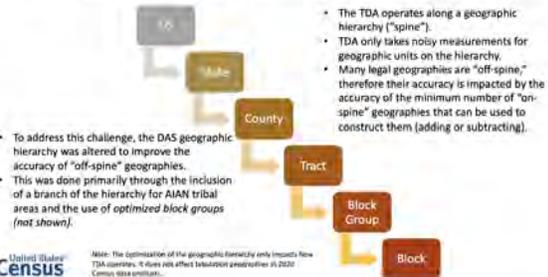


Multi-pass Post-processing

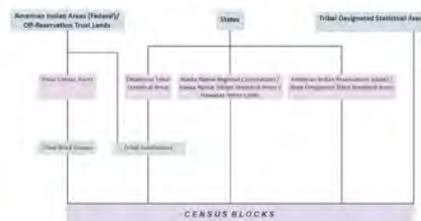
- The sparsity of many queries (i.e., prevalence of zeros and small counts) has the potential to introduce bias in TDA's post-processing
- To address the sparsity issue, TDA processing is performed in a series of passes
- At certain geographic levels, the algorithm constructs histograms for a subset of queries in a series of passes for that level, constraining the histogram for each pass to be consistent with the histogram produced in the prior pass
- Example for the P.L. 94-171 Redistricting Data Summary File:
 - Pass 1: Total Population
 - Pass 2: Remaining tabulations supporting P.L. 94-171 Redistricting Data

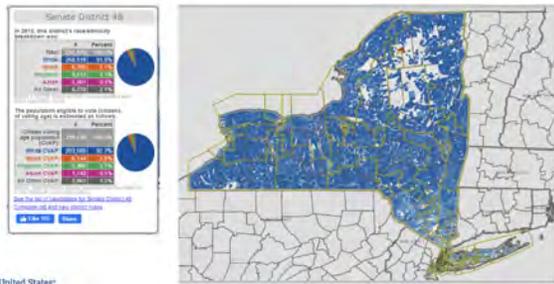


Tabulation Geographic Hierarchy



Hierarchy of American Indian, Alaska Native, and Native Hawaiian Areas





United States
Census

13



United States
Census

14

How to reconcile these statistics

- Construct error metrics of the form

$$Pr[|TDA - CEF| \leq \alpha] \geq 1 - \beta$$
- Less than α error with probability at least $1 - \beta$ for a target minimum population
- Statistical interpretation: absolute differences (=RMSE differences) greater than α are outside the $1 - \beta$ confidence interval
- A single statistic can be used to tune the redistricting application

$$\frac{\text{Population of Largest Race or Ethnic Group}}{\text{Total Population}}$$

- Calculated for the TopDown Algorithm (TDA) output and the 2020 Census (CEF)
- Implemented successfully for the production code release
- In the production data: minimum population of 200 to 249 for political areas and 450 to 499 for block groups to achieve 95% accuracy ($\alpha = 0.05$) at least 95% of the time ($\beta = 0.05$) See Wright and Imata (2021)

United States
Census

Figure 2. Mean Absolute Error of the County Total Population among the Least Populous Counties (Population Under 1,000) by Demonstration Data Product Vintage

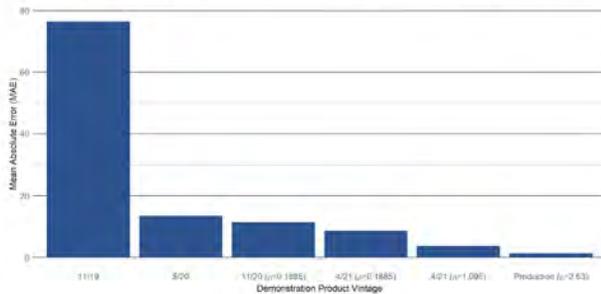
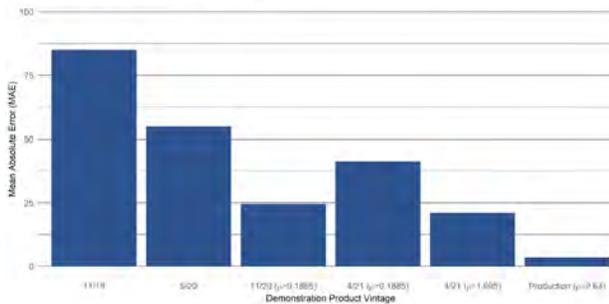


Figure 4. Mean Absolute Error of the Total Population among All Incorporated Places by Demonstration Data Product Vintage



What do the redistricting data do?

- Total differentially private measurements (queries): 16.6 billion
- Global $\rho = 2.63$ [(ϵ, δ) = (18.19, 10^{-10}) and infinitely many other pairs] U.S. persons and housing units
- Total block-level tables 29.4 million
- Total block-level statistics 3.4 billion
- Total independent block-level statistics 1.5 billion
- Accuracy of populations and largest race/ethnic group fit for redistricting and Voting Rights Act scrutiny for populations of at least 200-249, which is much smaller than legal entity subject to VRA

United States
Census

Figure 3. Mean Absolute Error of the Total Population for Federal American Indian Reservation/Off-Reservation Trust Lands by Demonstration Data Product Vintage

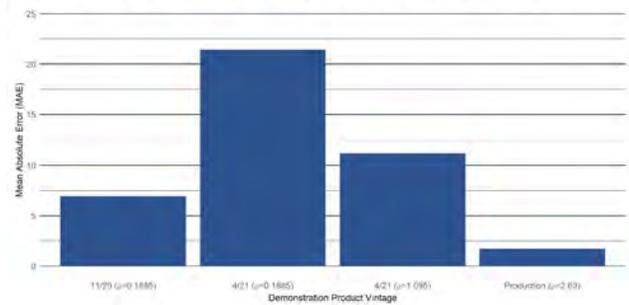
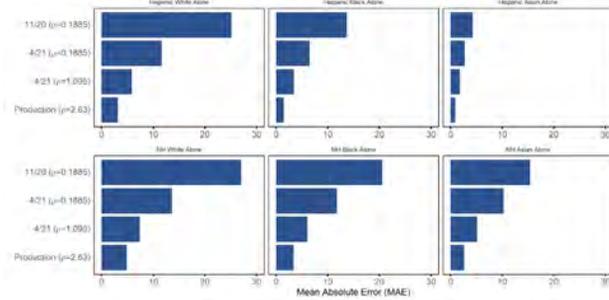


Figure 5. Mean Absolute Error of the Total Population among Tracts for Hispanic x Race Alone Populations by Demonstration Data Product Vintage



Addressing Other Biases

April 2021 PPMF		Production Settings	
Diversity Quintile	Mean Difference In Total Population	Diversity Quintile	Mean Difference In Total Population
0 - Least Diverse	5.04	0 - Least Diverse	-0.375
1	4.24	1	1.009
2	0.99	2	0.997
3	-2.21	3	-0.303
4 - Most Diverse	-8.07	4 - Most Diverse	-1.352

Block Groups		Tracts	
Diversity Quintile	Mean Difference In Total Population	Diversity Quintile	Mean Difference In Total Population
0 - Least Diverse	15.95	0 - Least Diverse	0.029
1	11.15	1	0.045
2	-3.01	2	0.000
3	-6.17	3	-0.020
4 - Most Diverse	-23.94	4 - Most Diverse	-0.053

Block-Level Inconsistencies Due to DAS-induced Uncertainty

Inconsistency	April 2021 Count of Blocks p=1.095	Production Settings Count of Blocks p=2.63
Occupied Housing Units > Household Population	203,519	303,984
Zero Occupied Housing Units; > 0 Household Population	674,598	505,840
Zero Household Population; > 0 Occupied Housing Units	77,947	148,836
Everyone in Block Under 18	90,534	163,884
> 10 Persons Per Household	87,342	121,376

Privacy-loss Budget Allocation (by geographic level)

Privacy-loss Budget Allocation 2021-06-08 Person Tables (Production Settings) United States		Privacy-loss Budget Allocation 2021-06-08 Units Tables (Production Settings) United States	
Global ρ	2.56	Global ρ	0.07
Global ϵ (incl. units)	18.19		
delta	10^{-11}		
	ρ Allocation by Geographic Level		ρ Allocation by Geographic Level
US	104/4099	US	1/205
State	1440/4099	State	1/205
County	447/4099	County	7/82
Tract	687/4099	Tract	364/1025
Optimized Block Group*	1256/4099	Optimized Block Group*	1759/4100
Block	165/4099	Block	99/820

Privacy-loss Budget Allocation (by query)

Query	Per Query ρ Allocation by Geographic Level					
	US	State	County	Tract	Optimized Block Group*	Block
TOTAL (1 cell)	3773/4097	1136/4097	1567/4102	1705/4099	5/4097	
CENRACE (63 cells)	52/4097	6/4097	10/4097	4/2051	3/4099	3/4097
HISPANIC (2 cells)	26/4097	6/4097	10/4097	5/4102	3/4099	5/4097
VOTINGAGE (2 cells)	26/4097	6/4097	10/4097	5/4102	3/4099	5/4097
FRINSTELEVELS (3 cells)	26/4097	6/4097	10/4097	5/4102	3/4099	5/4097
HWGQ (8 cells)	26/4097	6/4097	10/4097	5/4102	3/4099	5/4097
HISPANIC*CENRACE (126 cells)	130/4097	12/4097	28/4097	1933/4102	1055/4099	21/4097
VOTINGAGE*CENRACE (126 cells)	130/4097	12/4097	28/4097	10/2051	9/4099	21/4097
VOTINGAGE*HISPANIC (4 cells)	26/4097	6/4097	10/4097	5/4102	3/4099	5/4097
VOTINGAGE*HISPANIC*CENRACE (25 2 cells)	36/241	2/241	101/4097	67/4102	24/4099	71/4097
HWGQ*VOTINGAGE*						
HISPANIC*CENRACE (2,016 cells)	169/241	230/4097	754/4097	241/2051	1288/4099	3945/4097

Table 4 Accuracy of 2010 Census, enhanced Swap, and DP: mean absolute error (in persons) for age group population counts at the county level

Age group	2010 Census	Enhanced swap	DP
0-17 years	0	256.41	9.84
18-64 years	NA*	494.16	12.83
65 years and over	NA*	431.37	12.66

*Error statistics for the impact of swapping as applied to the published 2010 Census are confidential. The 2010 Census swapping algorithm kept the number of non-voting age individuals (0-17 years) invariant but did inject noise into the age groups within the voting age population. DRB clearance number CBDRB-FY22-DSEP-003. Data are from Devine & Spence (2022). Abbreviations: DP, differential privacy; DRB, Disclosure Review Board; NA, not available.

Table 5 Reidentification statistics for 2010 Census, enhanced swap, and DP

Reidentification Statistic	2010 Census	Enhanced swap	DP
Putative reidentification rate	97.0%	75.4%	44.4%
Confirmed reidentification rate	75.5%	46.6%	27.4%
Precision rate	77.8%	61.8%	61.7%
Precision for population uniques (nonmodal race)	81.4%	33.4%	21.0%

DRB clearance number CBDRB-FY22-DSEP-004. Data are from Abowd et al. (under review) released in Hawes (2022). External Matching File: Census Edited File. Abbreviations: DP, differential privacy; DRB, Disclosure Review Board.

• Tables 4 and 5 illustrate that TDA is a much more efficient disclosure avoidance mechanism for controlling accuracy and confidentiality than swapping with aggregation, as also shown in Abowd and Schmutte 2019.

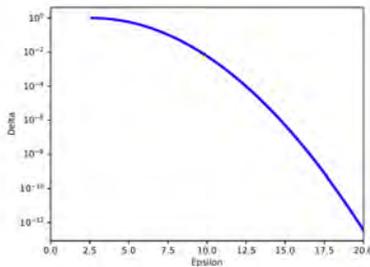


FIGURE 5. Bayesian (c, ϵ) curve under the semantics of Section 7.4.2 for $\rho = 2.63$

Source: Kifer et al. In preparation.

Significance Level	Power (Gaussian)	Power (DGM)	zCDP Upper Bound
0.01	0.032	0.032	0.037
0.05	0.12	0.12	0.14
0.10	0.21	0.21	0.24

TABLE 2. Block within Custom Block Group: Likelihood ratio test significance level/power tradeoff for block-level queries (1) if Gaussian noise is used, (2) if discrete Gaussian noise is used, (3) guaranteed upper bound if an arbitrary ρ -zCDP mechanism with $\rho = 0.1115$ is used.

$$\sup_{x>1} \frac{\alpha^x(1-\beta)^{1-x} + (1-\alpha)^x\beta^{1-x}}{e^{\rho x(1-x)}} \leq 1$$

where α is the level (probability of a Type I error), β is the probability of a Type II error, and $(1-\beta)$ is the power of the likelihood ratio test for correctly attaching a blockid to a record when block group, age, sex, race and ethnicity are known for zCDP. $H_0: N(0,1/2,\beta)$; $H_1: N(1,1/2,\beta)$

Source: Kifer et al. In preparation.

Figure B-13: Zachary Lipton - Adapting Predictors under Causally Structured Distribution Shift

Adapting Predictors under Causally Structured Distribution Shift

Zachary C. Lipton
 Assistant Professor
 Machine Learning Department
 & Operations Research
 Carnegie Mellon University

email: zlipton@cmu.edu | twitter: [@zacharylipton](https://twitter.com/zacharylipton)

Our lab's work in this line:

- Detecting & Correcting for Label Shift vs. Black Box Predictors (ICML 2018) (<https://arxiv.org/abs/1802.03916>)
- Falling Loudly (NeurIPS 2019) (<https://arxiv.org/abs/1810.11953>)
- What is the Effect of Importance Weighting in Deep Learning? (ICML 2019) (<https://arxiv.org/abs/1813.03372>)
- A Unified View of Label Shift Estimation (NeurIPS 2019) (<https://arxiv.org/abs/2003.07354>)
- Domain Adaptation with Asymmetrically-Relaxed Distribution Alignment (ICML 2019) (<https://arxiv.org/abs/1903.01689>)
- Learning Robust Representations by Penalizing Local Predictive Power (NeurIPS 2019) (<https://arxiv.org/abs/1905.13549>)
- Learning the Difference that Makes a Difference w/ Counterfactually Augmented Data (ICLR 2020) (<https://arxiv.org/abs/1909.12434>)
- Explaining the Efficacy of Counterfactually Augmented Data (ICLR 2021) (<https://arxiv.org/abs/2010.02114>)
- RATT: Leveraging Unlabeled Data to Guarantee Generalization (ICML 2021) (<https://arxiv.org/abs/2106.00303>)
- Mixture Proportion Estimation and PU Learning: A Modern Approach (NeurIPS 2021) (<https://arxiv.org/abs/2111.00980>)
- Online Estimation of Causal Effects by Deciding What to Observe (<https://arxiv.org/abs/2108.09265>)
- Unsupervised Learning Under Latent Label Shift (submitted) (<https://arxiv.org/abs/2002.13379.pdf>)
- Domain Adaptation under Open Set Label Shift (submitted) (<https://arxiv.org/pdf/2207.13048.pdf>)
- Domain Adaptation under Missingness Shift (forthcoming)

The Landscape of Distribution Shift

theoretically coherent work on idealized shift models

empirical deep learning efforts benchmark evaluation, heuristics

unpredictable shifts, limited faithfulness to any assumption

The Landscape of Distribution Shift

theoretically coherent work on idealized shift models

empirical deep learning efforts benchmark evaluation, heuristics

unpredictable shifts, limited faithfulness to any assumption

Structured Shifts

Fuzzy Shifts

The Landscape of Distribution Shift

theoretically coherent work on idealized shift models

empirical deep learning efforts benchmark evaluation, heuristics

unpredictable shifts, limited faithfulness to any assumption

Heuristics like matching representations or adapting batch statistics

GANIN 2016

NADO 2020

For problems, see: Domain Adaptation with Asymmetrically-Relaxed Distribution Alignment (ICML 2019) (<https://arxiv.org/abs/1903.01689>)

The Landscape of Distribution Shift

theoretically coherent work on idealized shift models

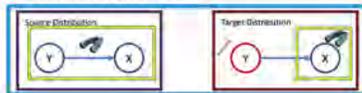
empirical deep learning efforts benchmark evaluation, heuristics

unpredictable shifts, limited faithfulness to any assumption

Koh & Sagawa 2020

Anatomy of a structured shift problem

- **Domains/Environments**—how many? how much data from each?
- **Structure**—model of data generating process
- **Visibility**—which variables are observed in each environment?
- **Manipulation rules**
 - What parts can/can't change?
 - By what amount?
 - In vs out-of-support?
- **Objective** (what to estimate)
- **Statistical Capabilities**
 - What relationships are estimable (& how well)?

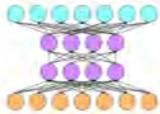


Some Examples of Structured Shift

- **Covariate Shift** — $P(Y|X)$ invariant, overlapping support $q(x) \subseteq p(x)$
- **Label Shift** — $P(X|Y)$ invariant, overlapping classes $q(y) \subseteq p(y)$
- **PU Learning** — $P(X|Y)$ invariant, + 1 new class: $P(Y=N) = 0, Q(Y=N) > 0$
- **Open Set Label Shift** — $P(Y|X)$ invariant, many prev classes, one new
- **Latent Label Shift** — $P(X|Y)$ invariant, many domains Q_i unsupervised
- **Missingness shift** — Source & target data missing at different rates

Two Obstacles to Practicality

- Identification is nice but we need practical estimators for high-dim data
- Graph-structure is not the only identification problem (eg, overlap violations)



- Assumptions too rigid, performance under fuzzy violations unknown

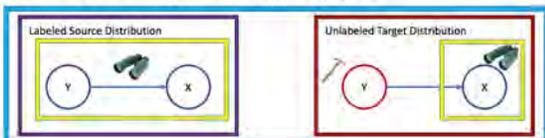
The Move: Leveraging Black Box Predictors

- No theory says we should be able to predict well (even on iid data) w high-dimensional, arbitrarily non-linear data (e.g. images, speech)
- However, we want to show that when it's possible to learn good iid classifiers, we can leverage these black boxes to get target classifiers

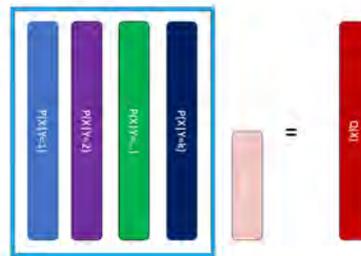


Label Shift (aka Target Shift)

- Assume $p(x,y)$ changes, but the conditional $p(x|y)$ is fixed
- $q(y, x) = q(y)p(x|y)$
- Corresponds to anticausal assumption, (y causes x)
- Assumptions: for all y such that $q(y) > 0, p(y) > 0$



Label Shift Identification



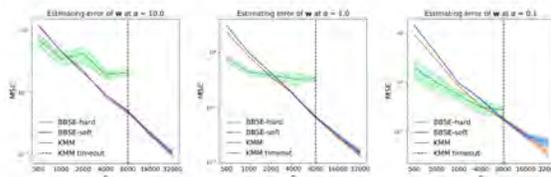
Approach 1: Black Box Shift Estimation (BBSE)

- Consistent estimator with interpretable error bounds
- Accuracy doesn't depend (directly) on ambient dimension
- Use black box for dim-reduction ($d \rightarrow 1$)
- Stronger $f \rightarrow$ tighter error bounds
- Inaccurate, uncalibrated $f \rightarrow$ BBSE still consistent

$$\hat{y} \quad C_{y|n} \quad y$$

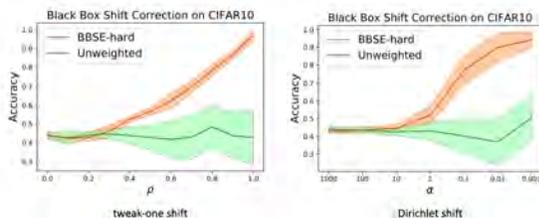
<https://arxiv.org/abs/1802.03916> (ZL, Wang, Smola—ICML 2018)

Estimation error in theory & practice



$$\|\hat{w} - w\|_2^2 \leq \frac{C}{\sigma_{\min}^2} \left(\frac{\|w\|^2 \log n}{n} + \frac{k \log m}{m} \right)$$

Black Box Shift Correction (CIFAR10 w IW-ERM)

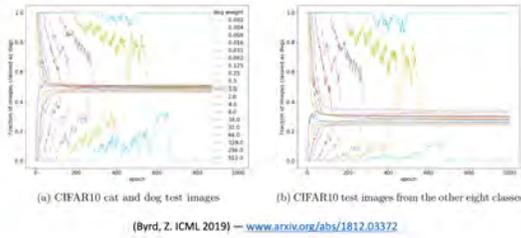


A General Pipeline for Detecting Shift

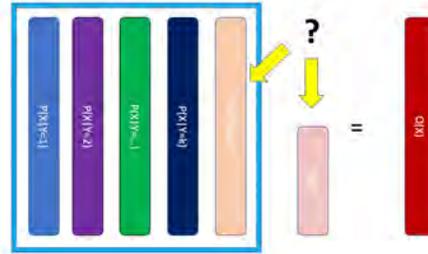


Failing Loudly: An Empirical Study of Methods for Detecting Dataset Shift (Rabanser, Gunneman, Z. — NeurIPS 2019) <https://arxiv.org/abs/1810.11053>

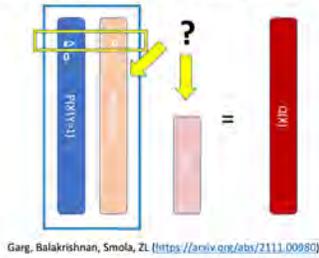
What is the effect of importance weighting in deep learning?



Label Shift with a New Class



PU Mixture Proportion Identification (irreducibility / positive subdomain)



MPE: Traditional Approaches

- Elkan and Nottó '08 discuss several approaches to estimate mixture proportion.
 - No guarantees and empirically perform bad
- Density estimation in input space [Ramaswamy+2013]. Curse of dimensionality in high dimensional datasets
- Recent methods that use classifier to reduce dimensionality
 - Need Bayes predictor for guarantees

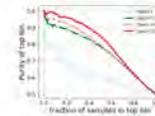
Estimation Strategy: Domain Discrimination

- Pool data
 - assign positive points label +,
 - assign unlabeled points label u
- Train a + vs. u classifier
- Assumption: when positive subdomain exists **most confidently predicted positive examples** are mostly positive



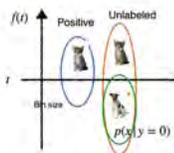
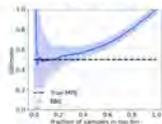
Our Estimator: Best Bin Estimation (BBE)

- Our approach: dimensionality reduction $f(x)$.
 - Use a classifier to transform $z = f(x)$
 - $q_u(z) = \alpha q_p(z) + (1 - \alpha)q_n(z)$
 - Estimate $\alpha^* = \min_z q_u(z)/q_p(z) \geq \alpha$
 - However, point estimates suffer from high variance
- Learn f to classify training v.s. test!
 - Pure top bin property: For $f(x) \geq t$, $P_n(f(x)) \approx 0$.
- Our estimator: $\hat{\alpha}(t) \approx P_u(f(x) \geq t) / P_p(f(x) \geq t)$



Identifying cutoff for the top bin

- How to choose t ?
 - A Naive approach: $t = 0.5$.
 - A better approach: $t = \arg \min_t \hat{\alpha}(t) + \text{conf}(\alpha(t))$
- Understanding Bias vs Variance tradeoff in the top bin



Theoretical Result

Theorem (Error rate of BBE)

Define $c^* = \arg \min_{c \in [0, 1]} q_u(c)/q_p(c)$. For $\min(n_p, n_u) \geq 2 \log(4/\delta)/q_p^*(c^*)$ and for every $\delta > 0$, the mixture proportion estimator $\hat{\alpha}$ (in Algo 1) satisfies the following with probability $1 - \delta$

$$|\hat{\alpha} - \alpha^*| \leq \frac{c_1}{q_p(c^*)} \left(\sqrt{\log(4/\delta)/n_u} + \sqrt{\log(4/\delta)/n_p} \right)$$

Obtaining the Classifier: Conditional Value Ignoring Risk (CVIR)

- We propose a simple objective

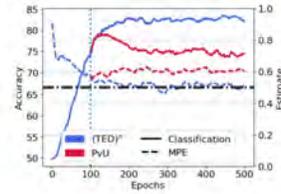
Algorithm 2 PU learning with Conditional Value Ignoring Risk (CVIR) objective

- Rank samples $x_n \in X_n$ according to their loss values $\ell(f_\theta(x_n), -1)$.
- $X_{n,1-\alpha} := X_n \setminus \alpha$ where $X_{n,1-\alpha}$ denote the lowest ranked $1-\alpha$ fraction of samples.
- Train model f_θ for one epoch on $(X_{n,1-\alpha}^1, X_n^1)$

- We can theoretically show that our loss will correctly discard positives from unlabeled when data is separable

Mixture Proportion Estimation and PU Learning: A Modern Approach

(TED)ⁿ: Combining CVIR and BBE



Mixture Proportion Estimation and PU Learning: A Modern Approach

MPE Results

Dataset	Model	(TED) ⁿ	BBE*	DeDup*	AlphaMax*	EN	KM2	TICE
Binarized CIFAR	ResNet	0.018	0.072	0.075	0.125	0.175		
	All Conv.	0.041	0.038	0.046	0.09	0.23	0.181	0.251
	FCN	0.184	0.175	0.151	0.3	0.355		
CIFAR Dog vs Cat	ResNet	0.074	0.120	0.113	0.17	0.205	0.11	0.203
	All Conv.	0.073	0.093	0.098	0.19	0.274		
Binarized MNIST	FCN	0.021	0.028	0.027	0.09	0.067	0.102	0.247
MNIST17	FCN	0.003	0.008	0.006	0.075	0.065	0.03	0.117
IMDb	BERT	0.008	0.011	0.016	0.07	0.12		

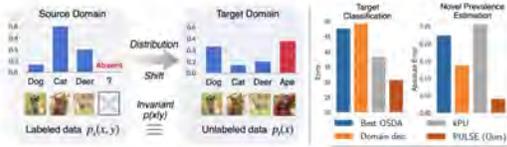
Mixture Proportion Estimation and PU Learning: A Modern Approach

Classification Results

Dataset	Model	(TED) ⁿ (unknown α)	CVIR (known α)	PvU* (known α)	DeDup* (unknown α)	mPU (known α)	uPU* (known α)	PvN
Binarized CIFAR	ResNet	82.7	82.6	78.3	78.4	76.8	75.8	86.9
	All Conv.	76.8	77.1	74.1	76.9	72.1	71.3	76.7
	FCN	63.2	65.9	61.4	62.5	63.9	64.8	65.1
CIFAR Dog vs Cat	ResNet	76.1	74.0	71.6	70.9	72.6	69.5	80.4
	All Conv.	72.2	71.0	70.1	70.5	68.4	65.2	77.9
Binarized MNIST	FCN	95.9	96.4	94.5	95.2	95.9	95.0	96.7
MNIST17	FCN	98.6	98.6	98.7	98.1	98.2	98.4	99.0
IMDb	BERT	87.6	87.4	86.1	87.3	86.2	85.9	89.1

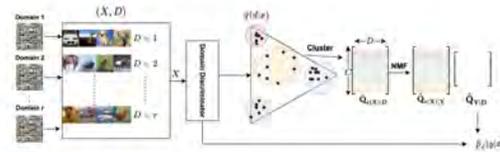
Mixture Proportion Estimation and PU Learning: A Modern Approach

Domain Adaptation under Open Set Label Shift (OSLS)



Garg, Balakrishnan, ZL (<https://arxiv.org/pdf/2207.13048.pdf>)

Unsupervised Learning w Latent Label Shift



Roberts, Maini, Garg, ZL (<https://arxiv.org/pdf/2207.13179.pdf>)

Thanks!

- Detecting & Correcting for Label Shift w Black Box Predictors (ICML 2018) (<https://arxiv.org/abs/1807.03910>)
- Falling Clouds (ICLR Deep ML 2019) (<https://arxiv.org/abs/1810.11853>)
- What is the Effect of Importance Weighting in Deep Learning? (ICML 2019) (<https://arxiv.org/abs/1812.03372>)
- A Unified View of Label Shift Estimation (NeurIPS 2020) (<https://arxiv.org/abs/2003.07568>)
- Domain Adaptation with Asymmetrically-Relaxed Distribution Alignment (ICML 2019) (<https://arxiv.org/abs/1903.01099>)
- Learning Robust Representations by Penalizing Local Predictive Power (NeurIPS 2019) (<https://arxiv.org/abs/1905.11549>)
- Learning the Difference that Makes a Difference w Counterfactually Augmented Data (ICLR 2020) (<https://arxiv.org/abs/1909.11435>)
- Explaining the Efficacy of Counterfactually Augmented Data (ICLR 2021) (<https://arxiv.org/abs/2010.02114>)
- RAT: Leveraging Unlabeled Data to Guarantee Generalization (ICML 2021) (<https://arxiv.org/abs/2105.09309>)
- Mixture Proportion Estimation And PU Learning: A Modern Approach (NeurIPS 2021) (<https://arxiv.org/abs/2111.02990>)
- Online Estimation of Causal Effects by Deciding What to Observe (<https://arxiv.org/abs/2108.09283>)
- Unsupervised Learning Under Latent Label Shift (<https://arxiv.org/pdf/2207.13048.pdf>)
- Domain Adaptation under Open Set Label Shift (<https://arxiv.org/pdf/2207.13048.pdf>)
- Domain Adaptation under Missingness Shift (forthcoming)

Figure B-14: Ruoxuan Xiong - Design and Analysis of Panel Data Experiments

Design and Analysis of Panel Data Experiments

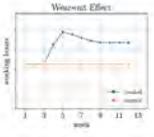
Ruoxuan Xiong
Emory University
Aug 17 @ AI/OR Workshop

[Optimal experimental design for staggered rollouts](#) with Susan Athey, Mohsen Bayati, Guido Imbens (XAB1'19)
[Bayesian methods for designing simultaneous temporal experiments](#) with Alex Chin, Sean Taylor, Susan Athey (NCTA'22)

Example: Marketplace experiments

A ridesharing platform tests an app feature to improve driver experience

- Conventional A/B testing suffers from network or contamination effect
- The impact of the feature may fade over time

Two key challenges

- Interference and network effects
- Carryover effects (instantaneous and lagged effects)

Panel data experiments can help addressing both challenges

Panel data experiments

- Run an experiment on N units (e.g. cities) for T time periods (e.g. days, weeks)
 - Repeated treatment decisions for the same N units
 - Repeated observations on the same N units

	day 1	day 2	day 3	day 4
NY	Control	Control	Treatment	Treatment
SF	Treatment	Treatment	Treatment	Treatment
ATL	Control	Treatment	Treatment	Treatment
⋮				

- Three main advantages
 - Increase sample size
 - Allow for heterogeneity in unit means and common time shocks
 - Study the impact of sustaining treatment over multiple periods

Treatment allocation schemes

Switchback experiments

	day 1	day 2	day 3	day 4
NY	Treatment	Treatment	Control	Control
SF	Treatment	Control	Control	Treatment
ATL	Control	Treatment	Control	Treatment
⋮				

Pro: Flexible
Con: Some practical constraints to switch back

Irreversible treatment adoption experiments

	day 1	day 2	day 3	day 4
NY	Control	Control	Treatment	Treatment
SF	Treatment	Treatment	Treatment	Treatment
ATL	Control	Treatment	Treatment	Treatment
⋮				

Simultaneous or staggered adoption pattern
Optimal design is a challenging problem!

Decision making problem

- Objective:** Most precisely estimate instantaneous and lagged effects
 - Can reduce sample size requirement and experimental cost!
- Decision variables:** Optimally choose the treatment times for each unit
 - Integer programming problem

	day 1	day 2	day 3	day 4	
NY	Control	Control	Treatment	Treatment	Z_{NY}^2
SF	Treatment	Treatment	Treatment	Treatment	Z_{SF}^2
ATL	Control	Treatment	Treatment	Treatment	Z_{ATL}^2
⋮					

Fixed-sample-size experiments

- N and T are set pre-experiment
- Consider a flexible outcome model
 - Nonstationarity and heterogeneity

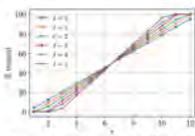
$$Y_{it}(Z) = \alpha_i + \beta_i + X_{it}' \cdot \theta_i + U_{it}' \cdot V_i + \tau_0 \cdot Z_{it} + \tau_1 \cdot Z_{it-1} + \dots + \tau_\ell \cdot Z_{it-\ell} + \epsilon_{it}$$

(Very) many fixed effects (unknown)
 Observational and/or common (with unknown/misspecified covariates)
 Binary treatment variables with unknown treatment effect parameters

- Formulate an integer program to solve Z that maximizes the estimation precision of instantaneous effect $\hat{\tau}_0$ and lagged effects $\hat{\tau}_1, \dots, \hat{\tau}_\ell$, post-experiment

Optimality conditions and choosing a design

- Provide the optimal conditions for the integer program [XAB1'19]
 - Optimal treated fraction at every time period
 - Covariate conditions
- Develop an algorithm to choose a design with provably guarantee to the optimal integer solution (within $1 + O(1/N^2)$) [XAB1'19]




Optimal treated fraction for a 12-period experiment

Sequential experiments

- N is fixed and T varies (we can early stop the experiment)
- More flexible and cost-effective!
- Key challenges
 - Experiment termination rule
 - What is an appropriate rule and how to implement the rule
 - Peeking challenge (Johari, Rosencor, Shalizi 2015)
 - Treatment effect estimation based by experiment termination rule
 - Infeasibility to optimize treatment times pre-experiment
 - Optimal solution depends on T



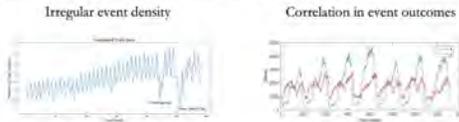
An initial stab at the problem

- We propose the Precision-Guided Adaptive Experimentation (P-GAE) algorithm (SABOTZ)
- Leverage ideas from dynamic programming, empirical Bayes, and sample splitting
- Key features
 - Adaptive treatment decisions
 - Precision-based experiment termination rule
 - Valid statistical inference post-experiment
- We provide theoretical guarantees for P-GAE
 - Asymptotic: consistency, normality, and efficiency



Micro-level perspective

- Micro-level data: raw data are events, like rider checking price, outcome is whether rider requested a ride
 - Large sample size, but analysis is more challenging (Johari, Li, Rosencor, Shalizi 2017)
- Additional considerations when analyzing event data



Additional considerations

- Additional considerations when analyzing event data



Error analysis and design of experiments

- Analyze the mean-squared error (MSE) in treatment effect estimation (SCT 3.2.2)
 - Bias affected by event density, spillover effects, simultaneous experiments
 - Variance affected by event density, correlation in event outcomes
- Study how partition time and space (irregularly) to minimize MSE (SCT 3.2.3)



Thank you!



Figure B-15: Yu Ding - Causal Inference in Engineering Applications

Causal Inference in Engineering Applications

Yu Ding, Ph.D.
Mike and Sugar Barnes Professor
Texas A&M University
Associate Director for Research Engagement
Texas A&M Institute of Data Science
Aug 17, 2022

Natural experiments and causal inference

- Natural experiments *"are observational studies and are not controlled in the traditional sense of a randomized experiment."* (source: Wikipedia)
- Causal inference aims at
 - determining which factors have a genuine cause-and-effect relationship with the response
 - quantifying the effect on the response due to the action taken.
- Huge impact in economics and social sciences.
- In some engineering applications, conducting controlled experiment is either too costly or infeasible. Then, causal inference becomes relevant.

Wind energy application

- Wind energy (the same for solar) is known as *variable* renewable energy.
- Its fuel input is not controllable. Nor are other environmental conditions, which affect the wind power production, too.

wind speed

wind direction

air density

wind dynamics (TI, shear)

other inflow conditions



wind power production

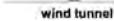
wind turbine power curve



Size matters

- Typical wind tunnels are an order of magnitude smaller than the size of commercial turbines.

A few meters in diameter



wind tunnel

Hub height: 80-100 m

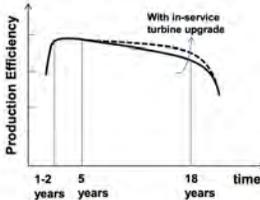


wind turbine (2 MW class)

- Controlled experiments can be done on lab-scaled small turbines, but extrapolating to the commercial setting causes huge inaccuracy.

Where does causal inference help?

- Both aims are relevant, but the second aim, i.e., the effect quantification, is particularly so.
- To counter a turbine's fast deterioration, a popular solution is to retrofit wind turbine, especially the blades. Such retrofit is known as **turbine upgrade**, meant to be **performance enhancement**.



Service providers sell all kinds of performance enhancement options.

<p>Blade Add-Ons</p> <ul style="list-style-type: none"> ✓ Vortex generators ✓ Gurney flap ✓ Leading edge protection 	<p>Software & Controller Upgrades</p> <ul style="list-style-type: none"> ✓ DEH software upgrades ✓ Parameter changes ✓ Controller upgrades
---	--

Stallard and Green (2014). "How does wind farm performance decline with age?" *Renewable Energy*, 66:775-786.

Tasker, Bechtel, Post. "Are Power Curve Upgrades Worth It? Measuring the ROI of Turbine Upgrades" *American Wind Energy Association Seminar*, Sept. 15, 2020.

Vortex generator installation



- Take vortex generators for example (picture courtesy of SMART BLADE® GmbH).
- Installing VGs on one turbine costs about \$10K, so for a 200-turbine farm, the total cost is \$2M.
- The key question is how soon can the wind farm operator pay back the expense?
- Depends on the increase in AEP after VG installation.

AEP = Annual Energy Production

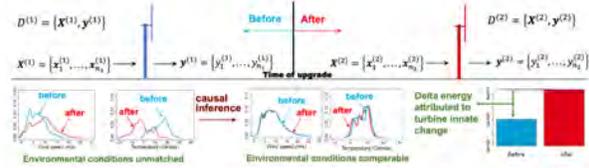


Causal inference for effect quantification

- If confirmed 1.5% – 3% improvement, as EDF advised, it makes a difference!
 3% improvement in AEP = \$1.2 million revenue increase
 (for a wind farm with 200 2MW-class turbines at 6 cents/kWh price)
- The general consensus is that the VG benefit in commercial settings is moderate, producing likely 1% – 5% increase in AEP.
- Defecting this small improvement and attributing its effect to VG installation is where causal inference can help with.

Covariate matching—A classical causal inference method

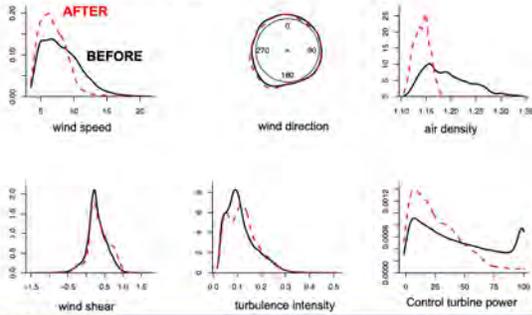
- Basic idea (Rubin 1973¹): carefully select data points for making the environmental conditions **probabilistically comparable** BEFORE and AFTER.
- Specific techniques (Shin et al. 2018¹¹):
 - 1 Hierarchical subgrouping
 - 2 Controlling for unmeasured factors
 - 3 One-to-one matching with replacement
 - 4 Robustness check on the order of matching



¹Rubin (1973). "Matching to remove bias in observational studies." *Biometrics*, 29: 159-183.
¹¹Shin, Ding, and Huang (2018). "Covariate matching methods for testing and quantifying wind turbine upgrades." *Annals of Applied Statistics*, 12: 1271-1292.

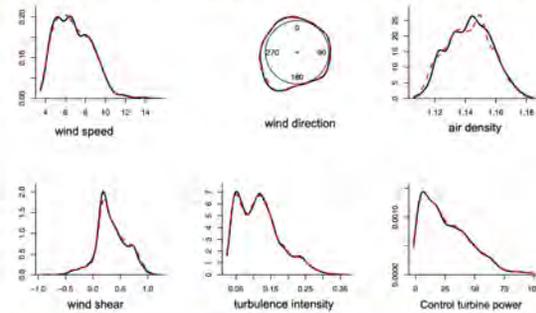
Matching results

- Two turbines are involved: a test turbine and a control turbine
- Before matching, the change in environmental conditions **confounds** the effect of a turbine upgrade



Matching results

- Two turbines are involved: a test turbine and a control turbine
- After matching, the density curves are probabilistically comparable.

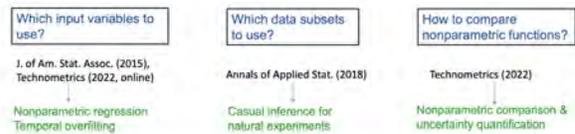


Further need for bias-reduction

- Even though a visual inspection shows reasonable alignment, there may still be a few percentage difference between two functional curves, large enough if compared with the anticipated upgrade effect.
- We explored a number of advanced options for covariate matching but those do not help and some even makes the bias worse.
- Bias reduction is still a pressingly needed.

Causal Inference Nested in a Three-Step Procedure

- Three methodological components.

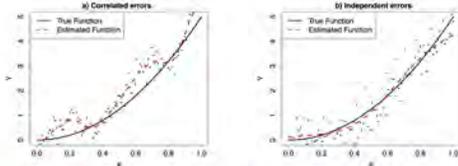


- The solution procedure[†] is implemented in both R and Python packages, named "DSWE".

[†]Ding, Kumar, Prakash, Ko, Liu, Liu, and Li (2021) "A case study of space-time performance comparison of wind turbines on a wind farm." *Renewable Energy*, 171: 735-746.

Step 1. Build a model that avoids temporal overfitting

- Cause-and-effect of certain variables is unmistakable, like wind speed, while conditional cause-and-effect of other variables are not so clear.
- When data are **autocorrelated**, deciding the conditional causal effect becomes tricky. Regular cross-validation tends to select an overfitting model.
- Our method, tempGP, is meant to extract the genuine effect despite the autocorrelation in data.



Prakash, Tu, and Ding (2022) "The temporal overfitting problem with applications in wind power curve modeling," *Technometrics*, online published.

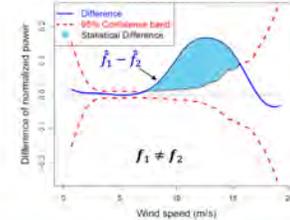
Step 3. Inference and UQ of difference in nonparametric functions

Let $f_1(\cdot)$ and $f_2(\cdot)$ be two nonparametric functions, representing the performances of two turbines of the same period, or two periods of the same turbine.

$$H_0 : f_1(x) = f_2(x) \quad \forall x \in \mathcal{X}$$

$$H_1 : f_1(x) \neq f_2(x) \quad \exists x \in \mathcal{X}$$

- 1 The functional estimates, $\hat{f}_1(\cdot)$ and $\hat{f}_2(\cdot)$, are used, as $f_1(\cdot)$ and $f_2(\cdot)$ are unknown.
- 2 Gaussian Process modeling provides foundation for uncertainty quantification.
- 3 Karhunen Loève (KL) expansion makes computation feasible.



Prakash, Tu, and Ding (2022) "Gaussian process aided function comparison using noisy scattered data," *Technometrics*, 64: 92-102.

A simulated study

The power in the test period is multiplied by a factor of $1 + r$ for those corresponding to wind speed greater than 9 m/s.

r	2%	3%	4%	5%	6%	7%	8%	9%
r^*	1.25%	1.87%	2.49%	3.11%	3.74%	4.36%	4.98%	5.60%
Δ in Ding et al. (2021)	1.12%	1.77%	2.73%	3.43%	3.69%	4.69%	5.16%	5.59%
Δ/r^*	0.90	0.95	1.10	1.10	0.99	1.08	1.04	1.00
UPG in Shin et al. (2018)	1.74%	2.21%	2.68%	3.16%	3.63%	4.11%	4.58%	5.05%
UPG/ r^*	1.39	1.18	1.08	1.02	0.97	0.94	0.92	0.90
DIFF in Lee et al. (2015)	1.97%	2.56%	3.15%	3.73%	4.30%	4.86%	5.42%	5.97%
DIFF/ r^*	1.58	1.37	1.27	1.20	1.15	1.11	1.09	1.07

- Ding et al. (2021): functional fitting + covariate matching + UQ.
- Shin et al. (2018): covariate matching only.
- Lee et al. (2015): functional fitting only.

Industrial case studies

- 2014, a blind study, with EDPR North America. [ASME TurboExpo Conference Proceedings, GT 2015, Montreal, Canada, June 15-19.]
- 2015-2016, with SMART BLADE[®] GmbH (Germany). [*Renewable Energy*, 2017, volume 113, pp. 1589–1597.]
- 2017-2018, with EDP Renewables (Spain/US). [A specialized R package named `gainML`].
The above studies focused on VG installation.
- 2019-2020, with Goldwind (China). Evaluation of software and control logic options. [*Renewable Energy*, 2021, volume 171, pp. 735–746].
- 2020-2021, with Wind Energy Institute of Canada (WEICan). Evaluation of leading-edge protection options on blades. [*Wind Energy*, 2022, volume 25, pp. 1203–1221].

Takeaways

- Causal inference to wind engineering: simple idea, real impact.
- For the aim of effect quantification, both bias reduction and uncertainty quantification are still challenging problems.
- Data characteristics deviating from standard iid assumptions raises new needs in causal inference.
- **Other applications.** Brian Denton recently published a perspective paper in *IJSE Transactions*, entitled "*Frontiers of medical decision-making in the modern age of data analytics*". The whole Section 2 is dedicated to the discussion of leveraging observational data to build IE/OR models.