
Beyond Ads: Sequential Decision-Making Algorithms in Public Policy

Peter Henderson*, Ben Chugg*, Brandon Anderson, Daniel E. Ho
Stanford University
phend@cs.stanford.edu, {benchugg,banderson,deho}@law.stanford.edu

Abstract

We explore the promises and challenges of employing sequential decision-making algorithms – such as bandits, reinforcement learning, and active learning – in the public sector. While such algorithms have been heavily studied in settings that are suitable for the private sector (e.g., online advertising), the public sector could greatly benefit from these approaches, but poses unique methodological challenges for machine learning. We highlight several applications of sequential decision-making algorithms in regulation and governance, and discuss areas for further research which would enable them to be more widely applicable, fair, and effective. In particular, ensuring that these systems learn rational, causal decision-making policies can be difficult and requires great care. We also note the potential risks of such deployments and urge caution when conducting work in this area. We hope our work inspires more investigation of public-sector sequential decision making applications, which provide unique challenges for machine learning researchers and can be socially beneficial.

1 Introduction

Sequential decision-making (SDM) algorithms come in many flavors, such as multi-armed bandits, reinforcement learning, and active learning. Together, these (overlapping) paradigms have been successfully applied to: content recommendation and ad placement [37, 55, 12], clinical trials [20, 5], robotics [50], and control of power systems [48, 52]. We argue that resource-constrained problems commonly faced by governments are a natural, but as of yet under-explored, application for SDM algorithms and that such public sector applications pose unique methodological challenges to the SDM literature. We highlight potential applications of SDM algorithms in the public sector,² and emphasize the social benefits of such deployments over the status quo.

The core of SDM is carefully balancing the explore-exploit trade-off. Should government agencies spend their resources to take advantage of the information they have now (“exploitation”), or spend some of those resources searching for better alternatives (“exploration”)? Too much of the former can result in missing opportunities and trends (e.g. discovering new forms of tax evasion or missing crucial evidence in a court case), while too much of the latter results in wasted resources. In many cases, as we will see, agencies already use what is effectively an SDM system without a formalization of these trade-offs. Yet formalization and improved methods could directly address the efficiency, transparency, and fairness of the process. We caution that in some cases SDM systems can be deployed in socially harmful ways; we discuss this further in Section 4 and urge careful consideration before deploying such a system.

*Equal contribution.

²We focus primarily on governance in the United States due to the expertise of authors, but most if not all of our examples have parallels in countries and regions around the world.

The remainder of the paper is organized as follows. Section 2 provides some brief background on the problem formulation. Section 3 then provides examples of problems faced by various agencies which are natural applications for SDM algorithms and highlights open methodological issues that deserve attention before widespread adoption. Section 4 addresses concerns regarding the use of these algorithms, and advice on when they should and should not be deployed.

2 Problem Formulation

The applications we discuss take place over discrete time steps or epochs $t = 1, \dots, T$. At each time t , we receive a set of N_t observations with features $\mathbf{X}_t = (x_{it})$, each of which has a hidden reward r_{it} . For example in the regulatory setting, reward will often be expressed as some quantitative measure of compliance (e.g., recovered tax payments or identification of environmentally non-compliant facilities). A *policy* decides which set of m_t observations to select and receive rewards for. Typically, $m_t \ll N_t$ though the set of observations and selection budget may change size over time ($N_t \neq N_{t+1}$ and $m_t \neq m_{t+1}$).

This framework is closest to the bandit formulation. In the multi-armed bandit setting, there are N “arms” to choose from, each associated with its own reward distribution. Every timestep, we choose an arm (or set of arms in the *batched* setting [45]). Our problem can be cast as a batched, *contextual bandit*, in which the context are the features for each observation. Alternatively, observations might be grouped in which case we can treat the problem as a contextual bandit in which arms are groups of observations [54]. If we additionally assume that rewards can be approximated as a functional form of observations, we can interpret the problem naturally as a structured bandit [41, 4].

We can also cast the active learning [51] problem in this framework. The goal of active learning is to maximize accuracy of the underlying model while making as few queries as possible. Given a pool of unlabelled observations (\mathbf{X}_t) the policy would decide which label (r_{it}) to reveal to maximize the model’s performance across the pool. Thus the reward is the reduction in generalization error from selecting a given arm. In some active learning formulations (like continuous active learning), the goal can instead be framed as successfully identifying all the labels of a given type [17]. In this case $r_{it} = 1$ if y_{it} is the goal label g_{it} .

We can adapt the reinforcement learning paradigm to this framework as well. For simplicity we define the reinforcement learning problem as an infinite-horizon discounted Markov Decision Process (MDP) [8, 47]. An MDP is defined as $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \beta, \gamma \rangle$ where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ is the reward function, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \mapsto \Delta(\mathcal{S})$ is the transition function, a kernel mapping state-action pairs to a probability distribution over $\mathcal{S} \times \mathcal{A}$, $\beta \in \Delta(\mathcal{S})$ is the initial state distribution, and $\gamma \in [0, 1)$ is the discount factor. In our case, more often than not, multiple actions will need to be taken per time-step and the transition function will not be known *a priori*. The state space will be constructed from \mathbf{X}_t , the set of arms available for selection. Furthermore, in many cases the true problem formulation may be better suited for the partially observable MDP framework [10], as the state space is unlikely to be fully observable. Unlike the bandit setting which maximizes myopic reward, the RL setting maximizes the expected discounted future return at each timestep $V^\pi(s) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t) | s_0 = s]$,

3 Applications & Obstacles

Here we discuss various applications of SDMs in public policy. Along with the examples, we highlight specific problems to be solved in each area to ensure that these algorithms be as effective, reliable, and fair as possible. For each example, we also give a possible formulation of the problem as an SDM algorithm.

This section should not be read as a comprehensive list of opportunities and challenges of applied ML. Much prior work has focused on challenges of deploying ML systems in various areas such as health, education, and the sciences (e.g., [24, 31, 7]). Privacy, fairness, and explainability, for instance, are important and well-documented concerns generally [44, 40, 2]. Here we focus on challenges that are of particular prominence for SDM in the public-sector.

3.1 Multi-Objective Decision-Making

Agencies can rarely afford to be single-minded. The Environmental Protection Agency (EPA), for example, might wish to identify the highest polluters and penalize them. But they may also need to accurately estimate the overall level of noncompliance by regulated facilities [42]. This is an oft encountered problem. In addition to targeting areas of risk, government agencies must be able to estimate total noncompliance in order to effectively guide policy.

We call this *population estimation*. Attempts to address multiple objectives in bandit and RL literature involve extending the reward scalar to a vector [19], but this discards powerful existing methods to evaluate specific peripheral objectives. Population estimation, for example, could be handled separately using survey design and sampling theory [38]. Integrating sampling design into the SDM process could leverage the problem structure for improved multi-objective performance.

Several open questions remain in this area. Can we simultaneously maximize reward and achieve unbiased population estimates? If not, can the trade-offs between these two objectives be made explicit? How can we add in tertiary objectives? Would such multi-objective designs benefit from a unified SDM framework which can combine research from active learning, bandits, reinforcement learning, and survey sampling?

Example 3.1.1 (Inspecting food safety at the FDA). The U.S. imports nearly 15 million yearly shipments of food.^a The Food Safety Modernization Act (FSMA) provides for inspections of imported food products and American-based food production facilities to ensure compliance with FDA standards.^b Understanding the number of at-risk facilities (the population total) is important for setting budgets and informing regulatory policy.

FORMULATION. The set of arms would be the facilities or shipments that can be audited, each with individual context \mathbf{X}_t . The policy learned would then select which arm to audit and be rewarded upon identification of non-compliance, secondarily estimating the total rate.

^a<https://www.fda.gov/media/120585/download>

^b<https://www.fda.gov/media/78021/download>

Example 3.1.2 (Technology Assisted Review during civil litigation). eDiscovery, short for electronic discovery, entails identifying relevant documents during the discovery process in legal proceedings. Since the body of potential evidence is often overwhelming, current state-of-the-art eDiscovery mechanisms use active learning to identify relevant documents [13, 17, 27]. We also note that identifying documents for FOIA requests uses a similar process in some cases.

FORMULATION. The policy sees a set of documents (arms), with X_t consisting of standard natural language processing features (e.g., TF-IDF vectors [30]). The policy selects a set of documents for lawyers to review and receives a reward if the document was responsive ($r_{it} = 1$) or non-responsive ($r_{it} = 0$). This is repeated across multiple rounds until there is confidence that all responsive documents have been identified. A random sampling mechanism is often employed to estimate the amount of remaining responsive documents [36] – the same as a secondary population estimation objective. Combining these two processes into one multi-objective system might help to reduce labeling costs.

3.2 Distribution Shifts in the Small Data Regime

The distributions of both policy-relevant features and their relationship with reward are rarely fixed in reality—behaviors adapt, policies change, exogenous economic shocks occur, etc. This *concept drift* emphasizes the exploration aspect of SDM, requiring careful selection just to keep up. Algorithms in public policy settings are doubly hampered as the amount of available samples is often small, potentially leaving little budget for reliably rewarding selections.

Recent work in both the bandit and active learning spaces focuses on either detecting concept drift or designing algorithms which are robust to it [11, 54, 32]. However, in public policy settings there are often *a priori* indicators of concept drift. Incorporating this structured drift information into

existing models could improve resilience and efficiency, particularly in small-data regimes. More research is needed to create reliable methods for incorporating such structural priors or external sources of information into SDM policies. Moreover, it is an open question as to which algorithms retain performance and fairness guarantees under such large discontinuous sources of drift.

Example 3.2.1 (Allocating public health resources during a pandemic). Infectious disease outbreaks can cause severe resource allocation issues. Effective distribution of tests and vaccines are of paramount importance for curtailing the spread of the disease. But conditions in this setting drift rapidly: one population could quickly become vaccinated and less vulnerable, or a new variant could quickly make another population more susceptible to infection. Could the incorporation of external information like trends from other regions help adjust to drift quickly?

FORMULATION. Recently, this problem has been formulated as a multi-armed bandit [15]. In this formulation, testing resources were allocated to different neighborhoods—each constituting an arm—rewarded by the number of COVID-19 positive individuals identified. Testing of incoming populations at the border has been similarly formulated as an RL problem, leading to more efficient detection of potentially infectious persons [6]. Successfully understanding the transition function in the RL setting could aid in handling distribution shift.

Example 3.2.2 (Optimizing tax audits at the IRS). Every year, the Internal Revenue Service (IRS) audits between 0.5% and 1.1% of individual taxpayer returns, out of a potential pool of hundreds of millions [49]. In addition to maximizing revenue, they must generate reliable population estimates (the “tax gap”), and conduct audits fairly [34]. At any given year a tax code change could affect the performance of a risk selection model.

FORMULATION. We might consider this a structured bandit problem. Each year the IRS selects a batch of taxpayers to audit from the population based on their featurized tax returns (\mathbf{X}_t). The reward received (r_t) is the difference between the taxpayer reported amount of taxes owed and the true adjusted amount after audit. r_t is related to taxpayer’s attributes x_{it} via some unknown function f and parameter θ_t : $r_{it} = f(x_{it}, \theta_t)$. To account for distribution shift the policy must estimate how much the underlying structure of the function f has changed given some update to the tax code.

3.3 Learning with and Identifying Corrupted Labels

In many public policy settings there is often a “human in the loop,” responsible for providing the labels or collecting the reward (e.g., auditors, inspectors). This can inject variance and subjectivity into the information received by the algorithm, which at best can make learning difficult and at worst can make an algorithmic approach infeasible (see Section 4). This leads to a general warning, illustrated by the examples below: if an algorithm is receiving information from heterogeneous and subjective sources, the variability in the information may swamp any true signal in the data and make machine learning a poor vehicle for solving the problem.

Example 3.3.1 (Prioritizing restaurant health inspections.) Recent work has suggested that machine and policy learning could aid in targeting restaurant health inspections [3, 25]. Like auditing at other agencies, this process could in theory be improved by SDM. But the restaurant health inspection process is highly stochastic. Recent evidence shows that two health inspectors could give drastically different scores to the same restaurant [28]. As a result, any algorithms using these noisy (or incorrect) labels will allocate more resources to areas where city health inspectors are most strict, perpetuating biases. Until SDM algorithms incorporate mechanisms to identify and correct for such noisy labels, such algorithms should not be deployed in the restaurant health inspection contexts.

FORMULATION. We consider each restaurant as an arm with features X_t indicating area, previous violation history, etc. The policy also takes into account the identity of the auditor

ξ_t . A reward r_t is given for identifying restaurants with key healthcode violations, but it is stochastic, heteroskedastic, and conditional on ξ_t . Some auditors ξ_t may provide incorrect labels with unknown frequency and unknown bias.

3.4 Feedback Loops

Ensign *et al.* [22] demonstrated that predictive policing algorithms were subject to runaway feedback loops: the algorithm would repeatedly focus on the same neighbourhoods. This is a danger in all SDM processes—a miscalibrated algorithm may focus only on the areas (or companies, individuals, etc.), at the expense of learning more about others. This effect has also been observed in health care and recommender systems [53, 1, 29]. Open research questions remain. How can you detect feedback loops and intervene? How can you determine when *not* to use an SDM system if there is a risk for such a feedback loop? Conversely, when does formalization of an existing ad hoc SDM system lead to a reduction of existing feedback loops?

Example 3.4.1 (Checking Environmental Compliance at the EPA). The Environmental Protection Agency (EPA) and state EPAs, must decide which facilities to inspect in order to ensure compliance with environmental laws such as the Clean Water Act, Clean Air Act, Safe Drinking Water Act, Toxic Substances Control Act, among others [14, 9]. Most of these inspections involve physical visitations which are time intensive and costly. As a result, only a handful of investigation and testing resources are allocated to various facilities. The same risk exists in this setting as in predictive policing—algorithms might fall into a local maximum in which they repeatedly allocate inspections to the same facilities.

FORMULATION. Each facility is an arm with features X_t indicating prior enforcement history, region, sensor metrics, calculated risk based on satellite imagery [14], etc. r_{it} is assigned for identifying high-polluting facilities.

Example 3.4.2 (Finding fraudulent claims at CMS). The Centers for Medicare and Medicaid Services (CMS) identifies and audits potentially fraudulent claims by providers.^a Currently, CMS uses the Fraud Prevention System (FPS), described as a “predictive analytics technology” and required by the Small Business Jobs Act of 2010. Such a system needs to avoid feedback loops in which providers with erroneous but harmless reports are continually targeted. As noted in the report, “[a]ny bias towards focusing on easily recoverable amounts could potentially skew program integrity efforts away from stopping some of the most egregious fraud.”

FORMULATION. Each provider is an arm with features X_t indicating prior enforcement history, degree of closeness to other providers based on social network (used in current systems, but potentially problematic for exacerbating feedback loops and unfairness), etc. r_{it} is assigned for successfully identifying fraudulent claims.

^ahttps://www.cms.gov/About-CMS/Components/CPI/Widgets/Fraud_Prevention_System_2ndYear.pdf

3.5 Causal Inference

An often-cited requirement of administrative decision-making is that it not be “arbitrary and capricious.” While it is unclear whether this standard would apply to algorithms adopted by agencies,³ the spirit of this rule is highly desirable for deployed SDM algorithms. For example, algorithms should not pick up on spurious correlations, or make arbitrary decisions. Often, when described in the context of algorithmic decision-making, explainability is considered a key requirement for government algorithms. But another desirable quality for preventing arbitrary decision-making is to ensure that systems take into account causal mechanisms. Many of the above challenges – for

³But see discussion in recent work on the matter [21, 16].

example learning from corrupted data or preventing feedback loops – can be thought of as ensuring a sampling regime that learns causal relationships, not simply correlations. Take, for example, the multi-objective nature of many public sector problems. To satisfy a reward maximizing objective, the sampling distribution of any SDM system will be biased, potentially interfering with learning causal relationships from unbalanced data. Without careful causal methods, this can also lead to feedback loops. Recent work has explored building causal mechanisms into SDM algorithms [35, 43, 33, 18, 56, 39]. But more work is needed to infer causal mechanisms in the face of challenges described above. Where little is known of the underlying causal graph, data is sparse, and causal structures can drift, SDM methods can quickly run into limits.

Example 3.5.1 (Identifying issues in judicial opinions at the SSA). The U.S. Social Security Administration (SSA) introduced the Insight natural language processing system to help analyze and spot potential errors in draft adjudicatory decisions [26]. Such systems may evolve into an SDM, where the system flags potentially problematic text in opinions and attorneys review the flagged text. One potential label is whether an appeals court reverses or remands the SSA decision. The outcome, however, may be confounded by legal representation, which exists only for a subset of cases and may matter substantially in whether an appeal is taken and succeeds. Instead of inferring spurious correlations, it will be important to identify deconfounded causal lexicons [46, 23] in the SDM setting.

FORMULATION. The policy must select from N opinions which are the most likely to be overturned on appeal and specifically which text is the problematic text. Lawyers review the document and identify potential errors, providing partial feedback to the model. Then several years later, the true reward is received on whether that opinion was overturned. In the meantime, the model must be updated with partial information.

4 Assessing and mitigating social harms

As evidenced by runaway feedback loops, poorly implemented SDM algorithms can have pernicious impacts. When should an SDM algorithm be deployed, and how should an agency weigh the costs and benefits of implementing such a system?

First, algorithmic impact assessments aim to foster reflection about the risks of adoption. In many instances, such reflection may rule out the deployment of algorithms with inherent harm (e.g., lethal weapons).

Second, in many instances, insufficient information exists at the time of the adoption decisions. There may be strong reasons to favor adoption, but benefits and costs are not precisely known. One benefit of SDMs is that they are sequentially implemented. Much like in the case of adaptive clinical trials – where sequential assignment enables researchers to limit the number of patients assigned to ineffective treatment arms – SDMs may enable ongoing assessments of an algorithmic impact. The stochastic component of an “explore” decision potentially enables researchers to ascertain the impact of an intervention, thereby curing the information deficit at the time of adoption. This epistemic benefit to SDMs may also enable more responsible adaptation, scaling, and discontinuation of algorithms in the public sector. But deployments must proceed with great caution.

5 Conclusions

We hope our work serves as a brief review of the potential for SDM in the public sector. Overall, bringing SDM algorithms to the sector has the potential to improve government services through efficiency, transparency, and fairness – if implemented with care. We argue that the research challenges in this sector are just as difficult, if not more so, than private sector challenges. And working on these public sector research challenges could bring large social benefits to many. But, as we repeat throughout, these applications must be handled with great care as they are often deployed in high-stakes settings. Results must be thoroughly vetted for robustness, fairness, and sound causal reasoning before deployment.

References

- [1] George Alexandru Adam, Chun-Hao Kingsley Chang, Benjamin Haibe-Kains, and Anna Goldenberg. Hidden risks of machine learning applied to healthcare: Unintended feedback loops between models and future data causing model degradation. In *Machine Learning for Healthcare Conference*, pages 710–731. PMLR, 2020.
- [2] Kasun Amarasinghe, Kit Rodolfa, Hemank Lamba, and Rayid Ghani. Explainable machine learning for public policy: Use cases, gaps, and research directions, 2021.
- [3] Susan Athey. Beyond prediction: Using big data for policy problems. *Science*, 355(6324):483–485, 2017.
- [4] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- [5] Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates. *Operations Research*, 68(1):276–294, 2020.
- [6] Hamsa Bastani, Kimon Drakopoulos, Vishal Gupta, Jon Vlachogiannis, Christos Hadjicristodoulou, Pagona Lagiou, Gkikas Magiorkinis, Dimitrios Paraskevis, and Sotirios Tsiodras. Efficient and targeted COVID-19 border testing via reinforcement learning, 2021.
- [7] Andrew L Beam and Isaac S Kohane. Big data and machine learning in health care. *Jama*, 319(13):1317–1318, 2018.
- [8] Richard Bellman. A markovian decision process. *Journal of mathematics and mechanics*, pages 679–684, 1957.
- [9] Elinor Benami, Reid Whitaker, Vincent La, Hongjin Lin, Brandon R. Anderson, and Daniel E. Ho. The distributive effects of risk prediction in environmental compliance: Algorithmic design, environmental justice, and public policy. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 90–105, New York, NY, USA, 2021. Association for Computing Machinery.
- [10] Anthony R Cassandra. A survey of POMDP applications. In *Working notes of AAAI 1998 fall symposium on planning with partially observable Markov decision processes*, volume 1724, 1998.
- [11] Emanuele Cavenaghi, Gabriele Sottocornola, Fabio Stella, and Markus Zanker. Non stationary multi-armed bandit: Empirical evaluation of a new concept drift-aware algorithm. *Entropy*, 23(3):380, 2021.
- [12] Wei Chen, Yajun Wang, and Yang Yuan. Combinatorial multi-armed bandit: General framework and applications. In *International Conference on Machine Learning*, pages 151–159. PMLR, 2013.
- [13] Rishi Chhatwal, Nathaniel Huber-Fliflet, Robert Keeling, Jianping Zhang, and Haozhen Zhao. Empirical evaluations of active learning strategies in legal document review. In *2017 IEEE International Conference on Big Data (Big Data)*, pages 1428–1437. IEEE, 2017.
- [14] Ben Chugg, Brandon Anderson, Seiji Eicher, Sandy Lee, and Daniel E Ho. Enhancing environmental enforcement with near real-time monitoring: Likelihood-based detection of structural expansion of intensive livestock farms. *International Journal of Applied Earth Observation and Geoinformation*, 103:102463, 2021.
- [15] Ben Chugg, Lisa Lu, Derek Ouyang, Benjamin Anderson, Raymond Ha, Alexis D’Agostino, Anandi Sujeer, Sarah L Rudman, Analilia Garcia, and Daniel E Ho. Evaluation of allocation schemes of COVID-19 testing resources in a community-based door-to-door testing program. In *JAMA Health Forum*, volume 2, pages e212260–e212260. American Medical Association, 2021.
- [16] Cary Coglianese and David Lehr. Transparency and algorithmic governance. *Admin. L. Rev.*, 71:1, 2019.

- [17] Gordon V Cormack and Maura R Grossman. Autonomy and reliability of continuous active learning for technology-assisted review. *arXiv preprint arXiv:1504.06868*, 2015.
- [18] Maria Dimakopoulou, Zhengyuan Zhou, Susan Athey, and Guido Imbens. Balanced linear contextual bandits. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):3445–3453, 2019.
- [19] Madalina M Drugan and Ann Nowe. Designing multi-objective multi-armed bandits algorithms: A study. In *The 2013 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2013.
- [20] Audrey Durand, Charis Achilleos, Demetris Iacovides, Katerina Strati, Georgios D Mitsis, and Joelle Pineau. Contextual bandits for adapting treatment in a mouse model of de novo carcinogenesis. In *Machine learning for healthcare conference*, pages 67–82. PMLR, 2018.
- [21] David Freeman Engstrom and Daniel E Ho. Algorithmic accountability in the administrative state. *Yale J. on Reg.*, 37:800, 2020.
- [22] Danielle Ensign, Sorelle A Friedler, Scott Neville, Carlos Scheidegger, and Suresh Venkatasubramanian. Runaway feedback loops in predictive policing. In *Conference on Fairness, Accountability and Transparency*, pages 160–171. PMLR, 2018.
- [23] Amir Feder, Katherine A Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E Roberts, et al. Causal inference in natural language processing: Estimation, prediction, interpretation and beyond. *arXiv preprint arXiv:2109.00725*, 2021.
- [24] Marzyeh Ghassemi, Tristan Naumann, Peter Schulam, Andrew L Beam, Irene Y Chen, and Rajesh Ranganath. A review of challenges and opportunities in machine learning for health. *AMIA Summits on Translational Science Proceedings*, 2020:191, 2020.
- [25] Edward L Glaeser, Andrew Hillis, Scott Duke Kominers, and Michael Luca. Crowdsourcing city government: Using tournaments to improve inspection accuracy. *American Economic Review*, 106(5):114–18, 2016.
- [26] Kurt Glaze, Daniel E. Ho, Gerald K. Ray, and Christine Tsang. Artificial intelligence for adjudication: The social security administration and ai governance. *Oxford Handbook on AI Governance*, 2022 (forthcoming).
- [27] Neel Guha, Peter Henderson, and Diego Zambrano. Vulnerabilities in discovery tech. *Harvard Journal of Law & Technology*, 2022 (forthcoming).
- [28] Daniel E Ho. Equity in the bureaucracy. *UC Irvine L. Rev.*, 7:401, 2017.
- [29] Ray Jiang, Silvia Chiappa, Tor Lattimore, András György, and Pushmeet Kohli. Degenerate feedback loops in recommender systems. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 383–390, 2019.
- [30] Daniel Jurafsky and James H Martin. Speech and language processing.
- [31] Anuj Karpatne, Imme Ebert-Uphoff, Sai Ravela, Hassan Ali Babaie, and Vipin Kumar. Machine learning for the geosciences: Challenges and opportunities. *IEEE Transactions on Knowledge and Data Engineering*, 31(8):1544–1554, 2018.
- [32] Bartosz Krawczyk and Alberto Cano. Adaptive ensemble active learning for drifting data stream mining. In *IJCAI*, pages 2763–2771, 2019.
- [33] Finnian Lattimore, Tor Lattimore, and Mark D Reid. Causal bandits: learning good interventions via causal inference. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 1189–1197, 2016.
- [34] Sarah B Lawsky. Fairly random: On compensating audited taxpayers. *Conn. L. Rev.*, 41:161, 2008.

- [35] Sanghack Lee and Elias Bareinboim. Structural causal bandits: where to intervene? *Advances in Neural Information Processing Systems* 31, 31, 2018.
- [36] Dan Li and Evangelos Kanoulas. When to stop reviewing in technology-assisted reviews: Sampling from an adaptive distribution to estimate residual relevant documents. *ACM Transactions on Information Systems (TOIS)*, 38(4):1–36, 2020.
- [37] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- [38] Sharon L Lohr. *Sampling: design and analysis*. Chapman and Hall/CRC, 2019.
- [39] Yangyi Lu, Amirhossein Meisami, and Ambuj Tewari. Causal bandits with unknown graph structure. *arXiv preprint arXiv:2106.02988*, 2021.
- [40] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021.
- [41] Adam J Mersereau, Paat Rusmevichientong, and John N Tsitsiklis. A structured multiarmed bandit problem and the greedy policy. *IEEE Transactions on Automatic Control*, 54(12):2787–2802, 2009.
- [42] U.S. Government Accountability Office. *Clean Water Act: EPA Needs to Better Assess and Disclose Quality of Compliance and Enforcement Data*. 2021.
- [43] Pedro A Ortega and Daniel A Braun. Generalized thompson sampling for sequential decision-making and causal inference. *Complex Adaptive Systems Modeling*, 2(1):1–23, 2014.
- [44] Nicolas Papernot, Patrick McDaniel, Arunesh Sinha, and Michael Wellman. Towards the science of security and privacy in machine learning. *arXiv preprint arXiv:1611.03814*, 2016.
- [45] Vianney Perchet, Philippe Rigollet, Sylvain Chassang, and Erik Snowberg. Batched bandit problems. *The Annals of Statistics*, 44(2):660–681, 2016.
- [46] Reid Pryzant, Kelly Shen, Dan Jurafsky, and Stefan Wagner. Deconfounded lexicon induction for interpretable social science. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1615–1625, 2018.
- [47] Martin L. Puterman. *Markov Decision Processes—Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, 1994.
- [48] Alberto Reyes, Matthijs TJ Spaan, and L Enrique Sucar. An intelligent assistant for power plants based on factored MDPs. In *2009 15th International Conference on Intelligent System Applications to Power Systems*, pages 1–6. IEEE, 2009.
- [49] Natasha Sarin and Lawrence H Summers. Shrinking the tax gap: approaches and revenue potential. Technical report, National Bureau of Economic Research, 2019.
- [50] Jonathan Scholz, Martin Levihn, Charles Isbell, and David Wingate. A physics-based model prior for object-oriented MDPs. In *International Conference on Machine Learning*, pages 1089–1097. PMLR, 2014.
- [51] Burr Settles. Active learning literature survey, 2009.
- [52] Mike Shann and Sven Seuken. Adaptive home heating under weather and price uncertainty using GPs and MDPs. In *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems*, pages 821–828, 2014.
- [53] Ayan Sinha, David F Gleich, and Karthik Ramani. Deconvolving feedback loops in recommender systems. *Advances in neural information processing systems*, 29:3243–3251, 2016.

- [54] Dennis Soemers, Tim Brys, Kurt Driessens, Mark Winands, and Ann Nowé. Adapting to concept drift in credit card transaction data streams using contextual bandits and decision trees. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [55] Liang Tang, Yexi Jiang, Lei Li, and Tao Li. Ensemble contextual bandits for personalized recommendation. In *Proceedings of the 8th ACM Conference on Recommender Systems*, pages 73–80, 2014.
- [56] Ruohan Zhan, Zhimei Ren, Susan Athey, and Zhengyuan Zhou. Policy learning with adaptively collected data. *arXiv preprint arXiv:2105.02344*, 2021.