

---

# Structural Causal Bandits: Where to Intervene?

---

**Sanghack Lee**

Department of Computer Science  
Purdue University  
lee2995@purdue.edu

**Elias Bareinboim**

Department of Computer Science  
Purdue University  
eb@purdue.edu

## Abstract

We study the problem of identifying the best action in a sequential decision-making setting when the reward distributions of the arms exhibit a non-trivial dependence structure, which is governed by the underlying causal model of the domain where the agent is deployed. In this setting, playing an arm corresponds to intervening on a set of variables and setting them to specific values. In this paper, we show that whenever the underlying causal model is not taken into account during the decision-making process, the standard strategies of simultaneously intervening on all variables or on all the subsets of the variables may, in general, lead to suboptimal policies, regardless of the number of interventions performed by the agent in the environment. We formally acknowledge this phenomenon and investigate structural properties implied by the underlying causal model, which lead to a complete characterization of the relationships between the arms' distributions. We leverage this characterization to build a new algorithm that takes as input a causal structure and finds a minimal, sound, and complete set of qualified arms that an agent should play to maximize its expected reward. We empirically demonstrate that the new strategy learns an optimal policy and leads to orders of magnitude faster convergence rates when compared with its causal-insensitive counterparts.

## 1 Introduction

The multi-armed bandit (MAB) problem is one of the prototypical settings studied in the sequential decision-making literature [Lai and Robbins, 1985, Even-Dar et al., 2006, Bubeck and Cesa-Bianchi, 2012]. An agent needs to decide which arm to pull and receives a corresponding reward at each time step while keeping the goal of maximizing its cumulative reward in the long run. The challenge is the inherent trade-off between exploiting known arms versus exploring new reward opportunities [Sutton and Barto, 1998, Szepesvári, 2010]. There is a wide range of assumptions underlying MABs, but in most of the traditional settings, the arms' rewards are assumed to be independent, which means that knowing the reward distribution of one arm has no implication to the reward of the other arms. Many strategies were developed to solve this problem, including classic algorithms such as  $\epsilon$ -greedy, variants of UCB [Auer et al., 2002, Cappé et al., 2013], and Thompson sampling [Thompson, 1933].

Recently, the existence of some non-trivial dependencies among arms has been acknowledged in the literature and studied under the rubric of *structured bandits*, which include settings such as linear [Dani et al., 2008], combinatorial [Cesa-Bianchi and Lugosi, 2012], unimodal [Combes and Proutiere, 2014], and Lipschitz [Magureanu et al., 2014], just to name a few. For example, a linear (or combinatorial) bandit imposes that an action  $x_t \in \mathbb{R}^d$  (or  $\{0, 1\}^d$ ) at a time step  $t$  incurs a cost  $\ell_t^\top x_t$ , where  $\ell_t$  is a loss vector chosen by, e.g., an adversary. In this case, an *index-based* MAB algorithm, oblivious to the structural properties, can be suboptimal.

In another line of investigation, rich environments with complex dependency structures are modeled explicitly through the use of causal graphs, where nodes represent decisions and outcome variables, and direct edges represent direct influence of one variable on another [Pearl, 2000]. Despite the

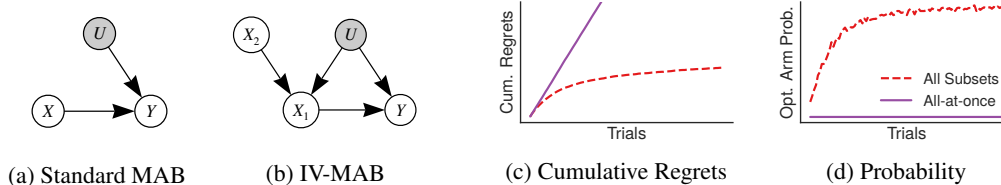


Figure 1: MAB problems as directed acyclic graphs where  $U$  is an unobserved variable. Plots of cumulative regrets and probability selecting an optimal arm when a MAB algorithm intervenes  $X_1$  and  $X_2$  simultaneously (All-at-once) or all subsets of  $\{X_1, X_2\}$  for IV-MAB. The IV-MAB is also used in the experimental section (see Appendix D [Lee and Bareinboim, 2018] for its parametrization).

apparent connection between MABs and causality, only recently has the use of causal reasoning been incorporated into the design of MAB algorithms. For instance, [Bareinboim et al., 2015] first explored the connection between causal models with unobserved confounders (UCs) and reinforcement learning, where latent factors affect both the reward distribution and the player’s intuition. The key observation used in the paper is that while standard MAB algorithms optimize based on the *do*-distribution (formally written as  $\mathbb{E}[Y|do(X)]$  or  $\mathbb{E}[Y_x]$ ), the simplest type of counterfactuals, this approach is dominated by another strategy using a more detailed counterfactual as the basis of the optimization process (i.e.,  $\mathbb{E}[Y_x|X = x']$ ); this general strategy was called regret decision criterion (RDC). This strategy was later extended to handle counterfactual distributions of higher dimensionality by [Forney et al., 2017]. Further, [Lattimore et al., 2016] and [Sen et al., 2017] studied the problem of best arm identification through importance weighting, where information on how playing arms influences the direct causes (parents, in causal terminology) of a reward variable is available. [Zhang and Bareinboim, 2017] leveraged causal graphs to solve the problem of off-policy evaluation in the presence of UCs. They noted that whenever UCs are present, traditional off-policy methods can be arbitrarily biased, leading to linear regret. They then showed how to solve the off-policy evaluation problem by incorporating the causal bounds into the decision-making procedure.<sup>1</sup> Overall, these works showed different aspects of the same phenomenon — whenever UCs are present in the real world, the expected guarantees provided by standard methods are no longer valid, which translates to an inability to converge to any reasonable policy. They then showed that convergence can be restored once the causal structure is acknowledged and used during the decision-making process.

In this paper, we focus on the challenge of identifying the best action in MABs where the arms correspond to interventions on an arbitrary causal graph, including when latent variables confound the observed relations (i.e., semi-Markovian causal models). To understand this challenge, we first note that a standard MAB can be seen as the simple causal model as shown in Fig. 1a, where  $X$  represents an arm (with  $K$  different values),  $Y$  the reward variable, and  $U$  the unobserved variable that generates the randomness of  $Y$ .<sup>2</sup> After a sufficiently large number of pulls of  $X$  (chosen by the specific algorithm),  $Y$ ’s average reward can be determined with high confidence.

Whenever a set of UCs affect more than one observed variable, however, novel, non-trivial challenges arise. To witness, consider the more involved MAB structure shown in Fig. 1b, where an unobserved confounder  $U$  affects both the action variable  $X_1$  and the reward  $Y$ . A naive approach for an algorithm to play such a bandit would be to pull arms in a combinatorial manner, i.e., combining both variables ( $X_1 \times X_2$ ) so that arms are  $D(X_1) \times D(X_2)$ , where  $D(X)$  is the domain of  $X$ . One may surmise that this is a valid strategy, albeit not the most efficient one. Somewhat unexpectedly, however, Fig. 1c shows that this is not the case — the optimal action comes from pulling  $X_2$  and ignoring  $X_1$ , while pulling  $\{X_1, X_2\}$  together would lead to subpar cumulative rewards (regardless of the number of iterations) since it simply cannot pull the optimal arm (Fig. 1d). After all, if one is oblivious to the causal structure and decides to take all intervenable variables as one (in this case,  $X_1 \times X_2$ ), indiscriminately, one may be doomed to learn a suboptimal policy.

<sup>1</sup>On another line of investigation, [Ortega and Braun, 2014] introduced a generalized version of Thompson sampling applied to the problem of adaptive control.

<sup>2</sup>In causal notation,  $Y \leftarrow f_Y(U, X)$ , which means that  $Y$ ’s value is determined by  $X$  and the realization of the latent variable  $U$ . If  $f_Y$  is linear, we would have a (stochastic) linear bandit. Our results do not constrain the types of structural functions, which is usually within nonparametric causal inference [Pearl, 2000, Ch. 7].

In this paper, we investigate this phenomenon, and more broadly, causal MABs with non-trivial dependency structure between the arms. More specifically, our contributions are as follows: (1) We formulate a SCM-MAB problem, which is a structured multi-armed bandit instance within the causal framework. We then derive the structural properties of a SCM-MAB, which are computable from any causal model, including arms' equivalence based on *do*-calculus [Pearl, 1995], and partial orderedness among sets of variables associated with arms in regards to the maximum rewards achievable. (2) We characterize a special set of variables called POMIS (possibly-optimal minimal intervention set), which is worth intervening based on the aforementioned partial orders. We then introduce an algorithm that identifies a complete set of POMISs so that only the subset of arms associated with them can be explored in a MAB algorithm. Simulations corroborate our findings.

**Big picture** The multi-armed bandit is a rich setting in which a huge number of variants has been studied in the literature. Different aspects of the decision-making process have been analyzed and well-understood in the last decades, which include different functional forms (e.g., linear, Lipschitz, Gaussian process), types of feedback experienced by the agent (bandit, semi-bandit, full), the adversarial or i.i.d. nature of the interactions, just to cite some of the most popular ones. Our study of SCM-MABs puts the causal dimension front and center in the map. In particular, we fully acknowledge the existence of a causal structure among the underlying variables (whenever not known *a priori*, see Footnote 3), and leverage the qualitative relations among them. This is in clear contrast with the prevailing practice that is more quantitative and, almost invariably, is oblivious to the underlying causal structure (as shown in Fig. 1a). We outline in Fig. 2 an initial map that shows the relationship between these dimensions; our goal here is not to be exhaustive, nor prescriptive, but to help to give some perspective. In this paper, we study bandits with no constraints over the underlying functional form (nonparametric, in causality language), i.i.d. stochastic rewards, and with an explicit causal structure acknowledged by the agent.

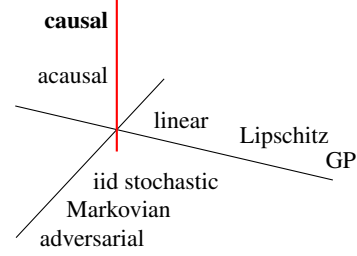


Figure 2: A bandit space with various dimensions (not all dimensions are shown)

### Preliminaries: notations and structural causal models

We follow the notation used in the causal inference literature. A capital letter is used for a variable or a mathematical object. The domain of  $X$  is denoted by  $D(X)$ . A bold capital letter is for a set of variables, e.g.,  $\mathbf{X} = \{X_i\}_{i=1}^n$ , while a lowercase letter  $x \in D(X)$  is a value assigned to  $X$ , and  $\mathbf{x} \in D(\mathbf{X}) = \times_{X \in \mathbf{X}} (D(X))$ . We denote by  $\mathbf{x}[\mathbf{W}]$ , values of  $\mathbf{x}$  corresponding to  $\mathbf{W} \cap \mathbf{X}$ . A graph  $G = \langle \mathbf{V}, \mathbf{E} \rangle$  is a pair of vertices  $\mathbf{V}$  and edges  $\mathbf{E}$ . We adopt family relationships — *pa*, *ch*, *an*, and *de* to denote parents, children, ancestors, and descendants of a given variable;  $Pa$ ,  $Ch$ ,  $An$ , and  $De$  extends *pa*, *ch*, *an*, and *de* by including the argument as the result, e.g.,  $Pa(X)_G = pa(X)_G \cup \{X\}$ . With a set of variables as argument,  $pa(\mathbf{X})_G = \bigcup_{X \in \mathbf{X}} pa(X)_G$  and similarly defined for other relations. We denote by  $\mathbf{V}(G)$  the set of variables in  $G$ .  $G[\mathbf{V}']$  for  $\mathbf{V}' \subseteq \mathbf{V}(G)$  is a vertex-induced subgraph where all edges among  $\mathbf{V}'$  are preserved. We define  $G \setminus \mathbf{X}$  as  $G[\mathbf{V}(G) \setminus \mathbf{X}]$  for  $\mathbf{X} \subseteq \mathbf{V}(G)$ .

We adopt the language of Structural Causal Models (SCM) [Pearl, 2000, Ch. 7]. An SCM  $M$  is a tuple  $\langle \mathbf{U}, \mathbf{V}, \mathbf{F}, P(\mathbf{U}) \rangle$ , where  $\mathbf{U}$  is a set of exogenous (unobserved or latent) variables and  $\mathbf{V}$  is a set of endogenous (observed) variables.  $\mathbf{F}$  is a set of deterministic functions  $\mathbf{F} = \{f_i\}$ , where  $f_i$  determines the value of  $V_i \in \mathbf{V}$  based on endogenous variables  $\mathbf{PA}_i \subseteq \mathbf{V} \setminus \{V_i\}$  and exogenous variables  $\mathbf{U}^i \subseteq \mathbf{U}$ , that is, e.g.,  $v_i \leftarrow f_i(\mathbf{pa}_i, \mathbf{u}^i)$ .  $P(\mathbf{U})$  is a joint distribution over the exogenous variables. A causal diagram  $G = \langle \mathbf{V}, \mathbf{E} \rangle$ , associated with  $M$ , is a tuple of vertices  $\mathbf{V}$  (the endogenous variables) and edges  $\mathbf{E}$ , where a directed edge  $V_i \rightarrow V_j \in \mathbf{E}$  if  $V_i \in \mathbf{PA}_j$ , and a bidirected edge between  $V_i$  and  $V_j$  if they share an unobserved confounder, i.e.,  $\mathbf{U}^i \cap \mathbf{U}^j \neq \emptyset$ . Note that  $pa(V_i)_G$  corresponds to  $\mathbf{PA}_i$ . Probability of  $Y = y$  when  $\mathbf{X}$  is held fixed at  $\mathbf{x}$  (i.e., intervened) is denoted by  $P(y|do(\mathbf{x}))$ , where intervention on  $\mathbf{X}$  is graphically represented by  $G_{\overline{\mathbf{X}}}$ , the graph  $G$  with incoming edges onto  $\mathbf{X}$  removed. We denote by  $CC(X)_G$  the *c-component* of  $G$  that contains  $X$  where a *c-component* is a maximal set of vertices connected with bidirected edges [Tian and Pearl, 2002]. We define  $CC(\mathbf{X})_G = \bigcup_{X \in \mathbf{X}} CC(X)_G$ . For a more detailed discussion on the properties of SCMs, we refer readers to [Pearl, 2000, Bareinboim and Pearl, 2016]. For all the proofs and appendices, please refer to the full technical report [Lee and Bareinboim, 2018].

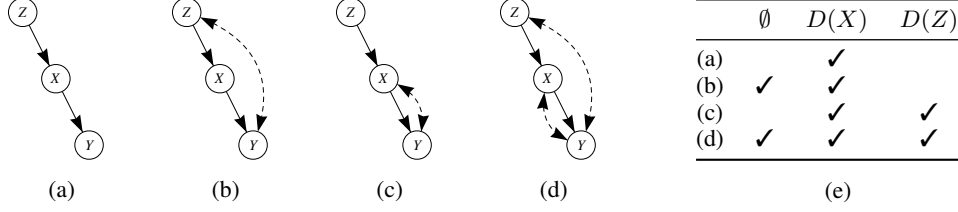


Figure 3: (a–d) Causal graphs such that  $\mu_x = \mu_{x,z}$ , and (e) non-dominated arms

## 2 Multi-armed bandits with structural causal models

We recall that MABs consider a sequential decision-making setting where pulling one of the  $K$  available arms at each round gives the player a stochastic reward from an unknown distribution associated with the corresponding arm. The goal is to minimize (maximize) the cumulative regret (reward) after  $T$  rounds. The mean reward of an arm  $a$  is denoted by  $\mu_a$  and the maximal reward is  $\mu^* = \max_{1 \leq a \leq K} \mu_a$ . We focus on the cumulative regret,  $\text{Reg}_T = T\mu^* - \sum_{t=1}^T \mathbb{E}[Y_{A_t}] = \sum_{a=1}^K \Delta_a \mathbb{E}[T_a(T)]$ , where  $A_t$  is the arm played at time  $t$ ,  $T_a(t)$  is the number of arm  $a$  has been played after  $t$  rounds, and  $\Delta_a = \mu^* - \mu_a$ .

We now can explicitly connect a MAB instance to its SCM counterpart. Let  $M$  be a SCM  $\langle \mathbf{U}, \mathbf{V}, \mathbf{F}, P(\mathbf{U}) \rangle$  and  $Y \in \mathbf{V}$  be a reward variable, where  $D(Y) \subseteq \mathbb{R}$ . The bandit contains arms  $\{\mathbf{x} \in D(\mathbf{X}) \mid \mathbf{X} \subseteq \mathbf{V} \setminus \{Y\}\}$ , a set of all possible interventions on endogenous variables except the reward variable. Each arm  $A_{\mathbf{x}}$  (or simply  $\mathbf{x}$ ) associates with a reward distribution  $P(Y|do(\mathbf{x}))$  where its mean reward  $\mu_{\mathbf{x}}$  is  $\mathbb{E}[Y|do(\mathbf{x})]$ . We call this setting a SCM-MAB, which is fully represented by the pair  $\langle M, Y \rangle$ . Throughout this paper, we assume that the causal graph  $G$  of  $M$  is fully accessible to the agent<sup>3</sup> although its parametrization is unknown: that is, an agent facing a SCM-MAB  $\langle M, Y \rangle$  plays arms with knowledge of  $G$  and  $Y$ , but not of  $\mathbf{F}$  and  $P(\mathbf{U})$ . For simplicity, we denote information provided to an agent playing a SCM-MAB by  $\llbracket G, Y \rrbracket$ . We now investigate some key structural properties that follow from the causal structure  $G$  of the SCM-MAB.

### Property 1. Equivalence among arms

We start by noting that *do*-calculus [Pearl, 1995] provides rules to evaluate invariances in the interventional space. In particular, we focus here on the Rule 3, which ascertains the condition such that a set of interventions does not have an effect on the outcome variable, i.e.,  $P(y|do(\mathbf{x}, \mathbf{z}), \mathbf{w}) = P(y|do(\mathbf{x}), \mathbf{w})$ . Since arms correspond to interventions (including the *null* intervention) and there is no contextual information, we consider examining  $P(y|do(\mathbf{x}, \mathbf{z})) = P(y|do(\mathbf{x}))$  through  $Y \perp\!\!\!\perp \mathbf{Z} \mid \mathbf{X}$  in  $G_{\overline{\mathbf{X}} \cup \mathbf{Z}}$ , which implies  $\mu_{\mathbf{x}, \mathbf{z}} = \mu_{\mathbf{x}}$ . If valid, this condition implies that it is sufficient to play only one arm among arms in the equivalence class.

**Definition 1** (Minimal Intervention Set (MIS)). A set of variables  $\mathbf{X} \subseteq \mathbf{V} \setminus \{Y\}$  is said to be a *minimal intervention set* relative to  $\llbracket G, Y \rrbracket$  if there is no  $\mathbf{X}' \subset \mathbf{X}$  such that  $\mu_{\mathbf{x}[\mathbf{X}']} = \mu_{\mathbf{x}}$  for every SCM conforming to the  $G$ .

For instance, the MISs corresponding to the causal graphs in Fig. 3 are  $\{\emptyset, \{X\}, \{Z\}\}$ , which do not include  $\{X, Z\}$  since  $\mu_x = \mu_{x,z}$ . The MISs are determined without considering the UCs in a causal graph. The empty set and all singletons in  $an(Y)_G$  are MISs for  $G$  with respect to  $Y$ . The task of finding the best arm among all possible arms can be reduced to a search within the MISs.

**Proposition 1** (Minimality). A set of variables  $\mathbf{X} \subseteq \mathbf{V} \setminus \{Y\}$  is a minimal intervention set for  $G$  with respect to  $Y$  if and only if  $\mathbf{X} \subseteq an(Y)_{G_{\overline{\mathbf{X}}}}$ .

All the MISs given  $\llbracket G, Y \rrbracket$  can be determined without explicitly enumerating  $2^{\mathbf{V} \setminus \{Y\}}$  while checking the condition in Prop. 1. We provide an efficient recursive algorithm enumerating the complete set of MISs given  $G$  and  $Y$  (Appendix A), which runs in  $O(mn^2)$  where  $m$  is the number of MISs.

<sup>3</sup>In settings where this is not the case, one can spend the first interactions with the environment to learn the causal graph  $G$  from observational [Spirtes et al., 2001] or experimental data [Kocaoglu et al., 2017].

## Property 2. Partial-orders among arms

We now explore the partial-orders among subsets of  $\mathbf{V} \setminus \{Y\}$  within the MISs. Given the causal diagram  $G$ , it is possible that intervening on some variables is *always* as good as intervening on another set of variables (regardless of the parametrization of the underlying model). Formally, there can be two different sets of variables  $\mathbf{W}, \mathbf{Z} \subseteq \mathbf{V} \setminus \{Y\}$  such that

$$\max_{\mathbf{w} \in D(\mathbf{W})} \mu_{\mathbf{w}} \leq \max_{\mathbf{z} \in D(\mathbf{Z})} \mu_{\mathbf{z}}$$

in every possible SCM conforming to  $G$ . If that is the case, it would be unnecessary (and possibly harmful in terms of sample efficiency) to play arms  $D(\mathbf{W})$ . We next define Possibly-Optimal MIS, which incorporates the partial-orderedness among subsets of  $\mathbf{V} \setminus \{Y\}$  into MIS denoting the optimal value for a  $\mathbf{X} \subseteq \mathbf{V} \setminus \{Y\}$  given a SCM by  $\mathbf{x}^*$ .

**Definition 2** (Possibly-Optimal Minimal Intervention Set (POMIS)). Given information  $\llbracket G, Y \rrbracket$ , let  $\mathbf{X}$  be a MIS. If there exists a SCM conforming to  $G$  such that  $\mu_{\mathbf{x}^*} > \forall \mathbf{Z} \in \mathbb{Z} \setminus \{\mathbf{x}\} \mu_{\mathbf{z}^*}$ , where  $\mathbb{Z}$  is the set of MISs with respect to  $G$  and  $Y$ , then  $\mathbf{X}$  is a *possibly-optimal minimal intervention set* with respect to the information  $\llbracket G, Y \rrbracket$ .

Intuitively, one may believe that the best action will be to intervene on the direct causes (parents) of the reward variable  $Y$ , since this would entail a higher degree of “controllability” of  $Y$  within the system. This, in fact, holds true if  $Y$  is not confounded with any of its ancestors, which includes the case where no unobserved confounders are present in the system (i.e., Markovian models).

**Proposition 2.** Given information  $\llbracket G, Y \rrbracket$ , if  $Y$  is not confounded with  $\text{an}(Y)_G$  via unobserved confounders, then  $\text{pa}(Y)_G$  is the only POMIS.

**Corollary 3** (Markovian POMIS). Given  $\llbracket G, Y \rrbracket$ , if  $G$  is Markovian, then  $\text{pa}(Y)_G$  is the only POMIS.

For instance, in Fig. 3a,  $\{\{X\}\}$  is the set of POMISs. Whenever unobserved confounders (UCs) are present<sup>4</sup>, on the other hand, the analysis becomes more involved. To witness, let us analyze the maximum achievable rewards of the MISs in the other causal diagrams in Fig. 3. We start with Fig. 3b and note that  $\mu_{z^*} \leq \mu_{x^*}$  since  $\mu_{z^*} = \sum_x \mu_x P(x|do(z^*)) \leq \sum_x \mu_x^* P(x|do(z^*)) = \mu_{x^*}$ . On the other hand,  $\mu_\emptyset$  is not comparable to  $\mu_{x^*}$ . For a concrete example, consider a SCM where the domains of variables are  $\{0, 1\}$ . Let  $U$  be the UC between  $Y$  and  $Z$  where  $P(U = 1) = 0.5$ . Let  $f_Z(u) = 1 - u$ ,  $f_X(z) = z$ , and  $f_Y(x, u) = x \oplus u$ , where  $\oplus$  is the exclusive-or function. If  $X$  is not intervened on,  $x$  will be  $1 - u$  yielding  $y = 1$  for both cases  $u = 0$  or  $u = 1$  so that  $\mu_\emptyset = 1$ . However, if  $X$  is intervened to either 0 or 1,  $y$  will be 1 only half the time since  $P(U = 1) = 0.5$ , which results in  $\mu_{x^*} = 0.5$ . We also provide in Appendix A a SCM such that  $\mu_\emptyset < \mu_{x^*}$  holds true. This model ( $\mu_\emptyset > \mu_{x^*}$ ) illustrates an interesting phenomenon — allowing an UC to affect  $Y$  freely may lead to a higher reward, which may be broken upon interventions. We now consider the different confounding structure shown in Fig. 3c (similar to Fig. 1b), where the variable  $Z$  lies outside of the influence of the UC associated with  $Y$ . In this case, intervening on  $Z$  leads to a higher reward,  $\mu_{z^*} \geq \mu_\emptyset$ . To witness, note that  $\mu_\emptyset = \sum_z \mathbb{E}[Y|z] P(z) = \sum_z \mu_z P(z) \leq \sum_z \mu_{z^*} P(z) = \mu_{z^*}$ . However,  $\mu_{z^*}$  and  $\mu_{x^*}$  are incomparable, which is shown through two models provided in Appendix A. Finally, we can add the confounders of the two previous models, which is shown in Fig. 3d. In this case, all three  $\mu_{x^*}$ ,  $\mu_{z^*}$ , and  $\mu_\emptyset$  are incomparable. One can imagine scenarios where the influence of the UCs are weak enough so that corresponding models produce results similar to Figs. 3a to 3c.

It’s clear that the interplay between the location of the intervened variable, the outcome variable, and the UCs entails non-trivial interactions and consequences in terms of the reward. The table in Fig. 3e highlights the arms that are contenders to generate the highest rewards in each model (i.e., each arm intervenes a POMIS to specific values), while intervening on a non-POMIS represents a waste of resources. Interestingly, the only parent of  $Y$ , i.e.,  $X$ , is not dominated by any other arms in any of the scenarios discussed. In words, this suggests that the intuition that controlling variables closer to  $Y$  is not entirely lost even when UCs are present; they are not the only POMIS, but certainly one of them. Given that more complex mechanisms cannot be, in general, ruled out, performing experiments would be required to identify the best arm. Still, the results of the table guarantee that the search can be refined so that MAB solvers can discard arms that cannot lead to profitable outcomes, and converge faster to playing the optimal arm.

<sup>4</sup>Recall that unobserved confounders are represented in the graph as bidirected dashed edges.



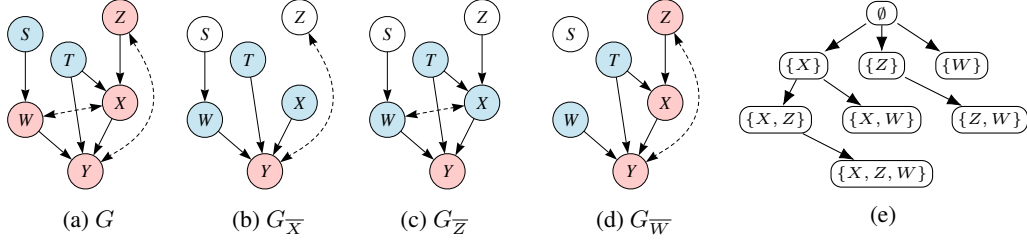


Figure 4: Causal graphs where pink and blue nodes are MUCT and IB, respectively. (Right most) A schematic showing an exploration order of subsets of variables.

### 3 Graphical characterization of POMIS

Our goal in this section is to graphically characterize POMISs. We will leverage the discussion in the previous section and note that UCs connected to a reward variable affect the reward distributions in a way that intervening on a variable outside the coverage of such UCs (including no UC) can be optimal — e.g.,  $\{X\}$  for Fig. 3a,  $\emptyset$  for Figs. 3b and 3d, and  $\{Z\}$  for Fig. 3c. We introduce two graphical concepts to help characterizing this property.

**Definition 3** (Unobserved-Confounders’ Territory). Given information  $\llbracket G, Y \rrbracket$ , let  $H$  be  $G[An(Y)_G]$ . A set of variables  $\mathbf{T} \subseteq \mathbf{V}(H)$  containing  $Y$  is called an *UC-territory* on  $G$  with respect to  $Y$  if  $De(\mathbf{T})_H = \mathbf{T}$  and  $CC(\mathbf{T})_H = \mathbf{T}$ .

An UC-territory  $\mathbf{T}$  is said to be *minimal* if no  $\mathbf{T}' \subset \mathbf{T}$  is an UC-territory. A minimal UC-Territory (MUCT) for  $G$  and  $Y$  can be constructed by extending a set of variables, starting from  $\{Y\}$ , alternatively updating the set with the c-component and descendants of the set.

**Definition 4** (Interventional Border). Let  $\mathbf{T}$  be a minimal UC-territory on  $G$  with respect to  $Y$ . Then,  $\mathbf{X} = pa(\mathbf{T})_G \setminus \mathbf{T}$  is called an *interventional border* for  $G$  with respect to  $Y$ .

The interventional border (IB) encompasses essentially the parents of the MUCT. For concreteness, consider Fig. 4a, and note that  $\{W, X, Y, Z\}$  is the MUCT for the causal graph with respect to  $Y$ , and the IB is  $\{S, T\}$  (marked in pink and blue in the graph, respectively). As its name suggests, MUCT is a set of endogenous variables governed by a set of UCs where at least one UC is adjacent to a reward variable. Specifically, the reward is determined by values of: (1) the UCs governing the MUCT; (2) a set of unobserved variables (other than the UCs) where each affects an endogenous variable in the MUCT; and (3) the IB. In other words, there is no UC interplaying *across* MUCT and its outside so that  $\mu_{\mathbf{x}} = \mathbb{E}[Y|\mathbf{x}]$  where  $\mathbf{x}$  is a value assigned to the IB  $\mathbf{X}$ . We now connect MUCT and IB with POMIS. Let  $MUCT(G, Y)$  and  $IB(G, Y)$  be, respectively, the MUCT and IB given  $\llbracket G, Y \rrbracket$ .

**Proposition 4.**  $IB(G, Y)$  is a POMIS given  $\llbracket G, Y \rrbracket$ .

The main strategy of the proof is to construct a SCM  $M$  where intervening on any variable in  $MUCT(G, Y)$  causes significant loss of reward. It seems that MUCT and IB can only identify a single POMIS given  $\llbracket G, Y \rrbracket$ . However, they, in fact, serve as basic units to identify all POMISs.

**Proposition 5.** Given  $\llbracket G, Y \rrbracket$ ,  $IB(G_{\overline{\mathbf{W}}}, Y)$  is a POMIS, for any  $\mathbf{W} \subseteq \mathbf{V} \setminus \{Y\}$ .

Prop. 5 generalizes Prop. 4 for when  $\mathbf{W} \neq \emptyset$  while taking care of UCs across  $MUCT(G_{\overline{\mathbf{W}}}, Y)$ , and its outside in the original causal graph  $G$ . See Fig. 4d for an instance, where  $IB(G_{\overline{\mathbf{W}}}, Y) = \{W, T\}$ . Intervening on  $W$  cuts the influence of  $S$  and the UC between  $W$  and  $X$ , while still allowing the UC to affect  $X$ .<sup>5</sup> Similarly, one can see in Fig. 4b that  $IB(G_{\overline{\mathbf{X}}}, Y) = \{T, W, X\}$  where intervening on  $X$  lets  $Y$  be the only element of MUCT making its parents an interventional border, hence, a POMIS. Note that  $pa(Y)_G$  is always a POMIS since  $MUCT(G_{\overline{pa(Y)_G}}, Y) = \{Y\}$  and  $IB(G_{\overline{pa(Y)_G}}, Y) = pa(Y)_G$ . With Prop. 5 one can enumerate the POMISs given  $\llbracket G, Y \rrbracket$  considering all subsets of  $\mathbf{V} \setminus \{Y\}$ . We show in the sequel that this strategy encompasses all the POMISs.

**Theorem 6.** Given  $\llbracket G, Y \rrbracket$ ,  $\mathbf{X} \subseteq \mathbf{V} \setminus \{Y\}$  is a POMIS if and only if  $IB(G_{\overline{\mathbf{X}}}, Y) = \mathbf{X}$ .

<sup>5</sup>Note that exogenous variables that do not affect more than one endogenous variable (i.e., non-UCs) are not explicitly represented in the graph.

---

**Algorithm 1** Algorithm enumerating all POMISs with  $\llbracket G, Y \rrbracket$ 

---

```
1: function POMISs( $G, Y$ )
2:    $\mathbf{T}, \mathbf{X} = \text{MUCT}(G, Y), \text{IB}(G, Y); H = G_{\overline{\mathbf{X}}}[\mathbf{T} \cup \mathbf{X}]$ 
3:   return  $\{\mathbf{X}\} \cup \text{subPOMISs}(H, Y, \text{reversed}(\text{topological-sort}(H)) \cap (\mathbf{T} \setminus \{Y\}), \emptyset)$ 
4: function SUBPOMISs( $G, Y, \pi, \mathbf{O}$ )
5:    $\mathbf{P} = \emptyset$ 
6:   for  $\pi_i \in \pi$  do
7:      $\mathbf{T}, \mathbf{X}, \pi', \mathbf{O}' = \text{MUCT}(G_{\overline{\pi_i}}, Y), \text{IB}(G_{\overline{\pi_i}}, Y), \pi^{i+1:|\pi|} \cap \mathbf{T}, \mathbf{O} \cup \pi^{1:i-1}$ 
8:     if  $\mathbf{X} \cap \mathbf{O}' = \emptyset$  then
9:        $\mathbf{P} = \mathbf{P} \cup \{\mathbf{X}\} \cup (\text{subPOMISs}(G_{\overline{\mathbf{X}}}[\mathbf{T} \cup \mathbf{X}], Y, \pi', \mathbf{O}') \text{ if } \pi' \neq \emptyset \text{ else } \emptyset)$ 
10:  return  $\mathbf{P}$ 
```

---

---

**Algorithm 2** POMIS-based kl-UCB

---

```
1: function POMIS-KL-UCB( $B, G, Y, f, T$ )
2:   Input:  $B$ , a SCM-MAB,  $G$ , a causal diagram;  $Y$ , a reward variable
3:    $\mathbf{A} = \bigcup_{\mathbf{X} \in \text{POMISs}(G, Y)} D(\mathbf{X})$ 
4:   kl-UCB( $B, \mathbf{A}, f, T$ )
```

---

Thm. 6 provides a graphical necessary and sufficient condition for a set of variables being a POMIS given  $\llbracket G, Y \rrbracket$ . This characterization allows one to determine all possible arms in a SCM-MAB that are worth intervening on, and, therefore, being free from pulling the other unnecessary arms.

## 4 Algorithmic characterization of POMIS

Although the graphical characterization provides a means to enumerate the complete set of POMISs given  $\llbracket G, Y \rrbracket$ , a naively implemented algorithm requires time exponential in  $|\mathbf{V}|$ . We construct an efficient algorithm (Alg. 1) that enumerates all the POMISs based on Props. 7 and 8 below and the graphical characterization introduced in the previous section (Thm. 6).

**Proposition 7.** Let  $\mathbf{T}$  and  $\mathbf{X}$  be the  $\text{MUCT}(G_{\overline{\mathbf{W}}}, Y)$  and  $\text{IB}(G_{\overline{\mathbf{W}}}, Y)$ , respectively, relative to  $G$  and  $Y$ . Then, for any  $\mathbf{Z} \subseteq \mathbf{V} \setminus \mathbf{T}$ ,  $\text{MUCT}(G_{\overline{\mathbf{X} \cup \mathbf{Z}}}, Y) = \mathbf{T}$  and  $\text{IB}(G_{\overline{\mathbf{X} \cup \mathbf{Z}}}, Y) = \mathbf{X}$ .

**Proposition 8.** Let  $H = G_{\overline{\mathbf{X}}}[\mathbf{T} \cup \mathbf{X}]$  where  $\mathbf{T}$  and  $\mathbf{X}$  are MUCT and IB given  $\llbracket G_{\overline{\mathbf{W}}}, Y \rrbracket$ , respectively. Then, for any  $\mathbf{W}' \subseteq \mathbf{T} \setminus \{Y\}$ ,  $H_{\overline{\mathbf{W}'}}$  and  $G_{\overline{\mathbf{W}' \cup \mathbf{W}}}$  yield the same MUCT and IB with respect to  $Y$ .

Prop. 7 allows one to avoid having to examine  $G_{\overline{\mathbf{W}}}$  for every  $\mathbf{W} \subseteq \mathbf{V} \setminus \{Y\}$ . Prop. 8 characterizes the recursive nature of MUCT and IB, where identification of POMISs can be evaluated by subgraphs. Based on these results, we design a recursive algorithm (Alg. 1) to explore subsets of  $\mathbf{V} \setminus \{Y\}$  with a certain order. See Fig. 4e for an example where subsets of  $\{X, Z, W\}$  are connected based on set inclusion relationship and an order of variables, e.g.,  $(X, Z, W)$ . That is, there exists a directed edge between two sets if (i) one set is larger than the other by a variable and (ii) the variable's index (as in the order) is larger than other variable's index in the smaller set. The diagram traces how the algorithm will explore the subsets following the edges, while effectively skipping nodes.

Given  $G$  and  $Y$ , POMISs (Alg. 1) computes a POMIS, i.e.,  $\text{IB}(G, Y)$ . Then, a recursive procedure subPOMISs is called with an order of variables (Line 3). Then subPOMISs examines POMISs by intervening on a single variable against the given graph (Line 6–9). If the IB ( $\mathbf{X}$  in Line 7) of such an intervened graph intersects with  $\mathbf{O}'$  (a set of variables that should be considered in other branch), then no subsequent call is made (Line 8). Otherwise, a subsequent subPOMISs call will take as arguments an MUCT-IB induced subgraph (Prop. 8), a refined order, and a set of variables not to be intervened in the given branch. For clarity, we provide a detailed working example in Appendix C with Fig. 4a where the algorithm explores only four intervened graphs ( $G, G_{\overline{\{X\}}}, G_{\overline{\{Z\}}}, G_{\overline{\{W\}}}$ ) and generates the complete set of POMISs  $\{\{S, T\}, \{T, W\}, \{T, W, X\}\}$ .

**Theorem 9** (Soundness and Completeness). *Given information  $\llbracket G, Y \rrbracket$ , the algorithm POMISs (Alg. 1) returns all, and only POMISs.*

The POMISs algorithm can be combined with a MAB algorithm, such as the kl-UCB, creating a simple yet effective SCM-MAB solver (see Alg. 2). kl-UCB satisfies  $\limsup_{n \rightarrow \infty} \frac{\mathbb{E}[\text{Reg}_n]}{\log(n)} \leq$

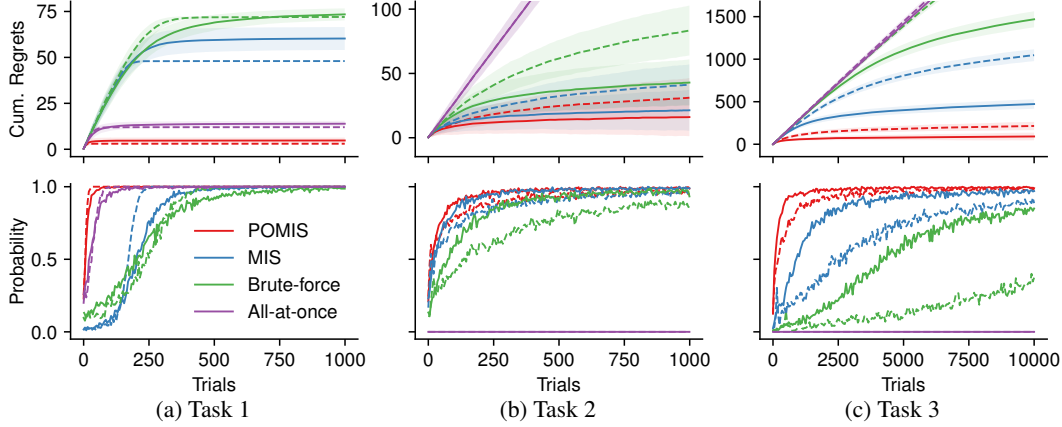


Figure 5: Comparisons across tasks (columns) with cumulative regrets (top) and optimal arm selection probability (bottom) with TS for solid and kl-UCB for dashed lines. Best viewed in color.

$\sum_{\mathbf{x}: \mu_{\mathbf{x}} < \mu^*} \frac{\mu^* - \mu_{\mathbf{x}}}{KL(\mu_{\mathbf{x}}, \mu^*)}$  where  $KL$  is Kullback-Leibler divergence between two Bernoulli distributions [Garivier and Cappé, 2011]. It is clear that the reduction in the size of arms will lower the upper bounds of the corresponding cumulative regrets.

## 5 Experiments

In this section, we present empirical results demonstrating that the selection of arms based on POMISs makes standard MAB solvers converge faster to an optimal arm. We employ two popular MAB solvers, kl-UCB, which enjoys cumulative regret growing logarithmically with the number of rounds [Cappé et al., 2013], and Thompson sampling (TS, [Thompson, 1933]), which has strong empirical performance [Kaufmann et al., 2012]. We considered four strategies for selecting arms, including POMISs, MISs, Brute-force, and All-at-once, where Brute-force evaluates all combinations of arms  $\bigcup_{\mathbf{x} \in \mathbf{V} \setminus \{Y\}} D(\mathbf{x})$ , and All-at-once considers intervening in all variables simultaneously,  $D(\mathbf{V} \setminus \{Y\})$ , oblivious to the causal structure and any knowledge about the action space. The performance of the eight ( $4 \times 2$ ) algorithms are evaluated relative to three different SCM-MAB instances (the detailed parametrizations are provided in Appendix D). We set the horizon large enough so as to observe near convergence, and repeat each simulation 300 times. We plot (i) the average cumulative regrets (CR) along with their respective standard deviations and (ii) the probability of an optimal arm being selected averaged over the repeated tests (OAP).<sup>6,7</sup>

**Task 1:** We start by analyzing a Markovian model. We note that by Cor. 3, searching for the arms within the parent set is sufficient in this case. The number of arms for POMISs, MISs, Brute-force, and All-at-once are 4, 49, 81, and 16, respectively. Note that there are 4 optimal arms within All-at-once arms — for instance, if the parent configuration is  $X_1 = x_1, X_2 = x_2$ , this strategy will also include combinations of  $Z_1 = z_1, Z_2 = z_2, \forall z_1, z_2$ . The simulated results are shown in Fig. 5a. CR at round 1000 with kl-UCB are 3.0, 48.0, 72, and 12 (in the order), and all strategies were able to find the optimal arms at this time. POMIS and All-at-once first reached 95% OAP at round 20 and 66, respectively. There are two interesting observations at this point. First, at an

<sup>6</sup>All the code is available at <https://github.com/sanghack81/SCMMAB-NIPS2018>

<sup>7</sup>One may surmise that combinatorial bandit (CB) algorithms can be used to solve SCM-MAB instances by noting that an intervention can be encoded as a binary vector, where each dimension in the vector corresponds to intervening on a single variable with a specific value. However, the two settings invoke a very different set of assumptions, which makes their solvers somewhat difficult to compare in some reasonably fair way. For instance, the current generation of CB algorithms is oblivious to the underlying causal structure, which makes them resemble very closely the Brute-force strategy, the worst possible method for SCM-MABs. Further, the assumption of linearity is arguably one of the most popular considered by CB solvers. The corresponding algorithms, however, will be unable to learn the arms' rewards properly since a SCM-MAB is nonparametric, making no assumption about the underlying structural mechanisms. These are just a few immediate examples of the mismatches between the current generation of algorithms for both causal and combinatorial bandits.



early stage, OAP for MISs is smaller than Brute-force since it has only 1 optimal arm among 49 arms, while Brute-force has 9 among 81. The advantage of employing MIS over Brute-force is only observed after a sufficiently large number of plays. More interestingly, POMIS and All-at-once both have the common optimal to non-optimal arms-ratio (1:3 versus 4:12), however, POMIS dominates All-at-once since the agent can learn better about the mean reward of the optimal arm while playing non-optimal arms less. Naturally, this translates into less variability and additional certainty about the optimal arm even in Markovian settings.

**Task 2:** We consider the setting known as instrumental variable (IV), which was shown in Fig. 3c. The optimal arm in this simulation is setting  $Z = 0$ . The number of arms for the four strategies is 4, 5, 9, and 4, respectively. The results are shown in Fig. 5b. Since the All-at-once strategy only considers non-optimal arms (i.e., pulling  $Z, X$  together), it incurs in a linear regret without selecting an optimal arm (0%). CR (and OAP) at round 1000 with TS are POMIS 16.1 (98.67%), MIS 21.4 (99.00%), Brute-force 42.9 (93.33%), and All-at-once 272.1 (0%). At round 5000, where Brute-force nearly converged, the ratio of CRs for POMIS and Brute-force is  $\frac{54.2}{18.1} = 2.99 \approx 2.67 = \frac{9-1}{4-1}$ . POMIS, MIS, and Brute-force first hits 95% OAP at 172, 214, and 435.

**Task 3:** Finally, we study the more involved scenario shown in Fig. 4a. In this case, the optimal arm is intervening on  $\{S, T\}$ , which means that the system should follow its natural flow of UCs, which All-at-once is unable to “pull.” There are 16, 75, 243, and 32 arms for the strategies (in the order). The results are shown in Fig. 5c. The CR (and OAP) at round 10000 with TS are POMIS 91.4 (99.0%), MIS 472.4 (97.0%), Brute-force 1469.0 (85.0%), and All-at-once 2784.8 (0%). Similarly, the ratio (in round 10000) is  $\frac{1469.0}{91.4} = 16.07 \approx 16.13 = \frac{243-1}{16-1}$  which is expected to increase since Brute-force is not yet converged at the moment. Only POMIS and MIS achieved OAP of 95% first in 684 and 3544 steps, respectively.

We start by noticing that the reduction in the CRs is approximately proportional to the reduction in the number of non-optimal arms pulled by (PO)MIS by the corresponding algorithm, which makes the POMIS-based solver the clear winner throughout the simulations. It’s still not inconceivable that the number of arms examined by All-at-once is smaller than for POMIS in a specific SCM-MAB instance, which would entail a lower CR to the former. However, such a lower CR in some instances does not constitute any sort of assurance since arms excluded from All-at-once, but included in POMIS, can be optimal in some SCM-MAB instance conforming to  $\llbracket G, Y \rrbracket$ . Furthermore, a POMIS-based strategy always dominates the corresponding MIS and Brute-force ones. These observations together suggest that, in practice, a POMIS-based strategy should be preferred given that it will always converge and will usually be faster than its counterparts. Remarkably, there is an interesting trade-off between having knowledge of the causal structure versus not knowing the corresponding dependency structure among arms, and potentially incurring in linear regret (All-at-once) or exponential slow-down (Brute-force). In practice, for the cases in which the causal structure is unknown, the pull of the arms themselves can be used as experiments and could be coupled with efficient strategies to simultaneously learn the causal structure [Kocaoglu et al., 2017].

## 6 Conclusions

We studied the problem of deciding whether an agent should perform a causal intervention and, if so, which variables it should intervene upon. The problem was formalized using the logic of structural causal models (SCMs) and formalized through a new type of multi-armed bandit called SCM-MABs. We started by noting that whenever the agent cannot measure all the variables in the environment (i.e., unobserved confounders exist), standard MAB algorithms that are oblivious to the underlying causal structure may not converge, regardless of the number of interventions performed in the environment. (We note that the causal structure can easily be learned in a typical MAB setting since the agent always has interventional capabilities.) We introduced a novel decision-making strategy based on properties following the *do*-calculus, which allowed the removal of redundant arms, and the partial-orders among the sets of variables existent in the underlying causal system, which led to the understanding of the maximum achievable reward of each interventional set. Leveraging this new strategy based on the possibly-optimal minimal intervention sets (called POMIS), we developed an algorithm that decides whether (and if so, where) interventions should be performed in the underlying system. Finally, we showed by simulations that this causally-sensible strategy performs more efficiently and more robustly than their non-causal counterparts. We hope that formal machinery and the algorithms developed here can help decision-makers to make more principled and efficient decisions.

## Acknowledgments

This research is supported in parts by grants from IBM Research, Adobe Research, NSF IIS-1704352, and IIS-1750807 (CAREER).

## References

- Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2/3):235–256, 2002.
- Elias Bareinboim and Judea Pearl. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352, 2016.
- Elias Bareinboim, Andrew Forney, and Judea Pearl. Bandits with unobserved confounders: A causal approach. In *Advances in Neural Information Processing Systems 28*, pages 1342–1350. 2015.
- Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Maillard, Rémi Munos, and Gilles Stoltz. Kullback-Leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, 41(3):1516–1541, 2013.
- Nicolò Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404 – 1422, 2012.
- Richard Combes and Alexandre Proutiere. Unimodal bandits: Regret lower bounds and optimal algorithms. In *International Conference on Machine Learning*, pages 521–529, 2014.
- Varsha Dani, Thomas P. Hayes, and Sham M. Kakade. Stochastic linear optimization under bandit feedback. In *Proceedings of Conference On Learning Theory (COLT)*, pages 355–366, 2008.
- Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of machine learning research*, 7(Jun):1079–1105, 2006.
- Andrew Forney, Judea Pearl, and Elias Bareinboim. Counterfactual data-fusion for online reinforcement learners. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1156–1164, 2017.
- Aurélien Garivier and Olivier Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *Proceedings of the 24th annual Conference On Learning Theory*, pages 359–376, 2011.
- Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. Thompson sampling: An asymptotically optimal finite-time analysis. In *Algorithmic Learning Theory*, pages 199–213, 2012.
- Murat Kocaoglu, Karthikeyan Shanmugam, and Elias Bareinboim. Experimental design for learning causal graphs with latent variables. In *Advances in Neural Information Processing Systems 30*, pages 7021–7031, 2017.
- Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- Finnian Lattimore, Tor Lattimore, and Mark D. Reid. Causal bandits: Learning good interventions via causal inference. In *Advances in Neural Information Processing Systems 29*, pages 1181–1189. 2016.
- Sanghack Lee and Elias Bareinboim. Structural causal bandits: Where to intervene? Technical Report R-36, Purdue AI Lab, Department of Computer Science, Purdue University, 2018.
- Stefan Magureanu, Richard Combes, and Alexandre Proutiere. Lipschitz bandits: Regret lower bound and optimal algorithms. In *Proceedings of The 27th Conference on Learning Theory*, pages 975–999, 2014.

- Pedro A. Ortega and Daniel A. Braun. Generalized Thompson sampling for sequential decision-making and causal inference. *Complex Adaptive Systems Modeling*, 2(2), 2014.
- Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995.
- Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, 2000. Second ed., 2009.
- Rajat Sen, Karthikeyan Shanmugam, Alexandros G. Dimakis, and Sanjay Shakkottai. Identifying best interventions through online importance sampling. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3057–3066, 2017.
- Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, Prediction, and Search*. A Bradford Book, 2001.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 1998.
- Csaba Szepesvári. *Algorithms for reinforcement learning*. Morgan and Claypool, 2010.
- William R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- Jin Tian and Judea Pearl. A general identification condition for causal effects. In *Proceedings of the Eighteenth National Conference on Artificial Intelligence*, pages 567–573, 2002.
- Junzhe Zhang and Elias Bareinboim. Transfer learning in multi-armed bandits: A causal approach. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pages 1340–1346, 2017.