

---

# Casual Bandits

---

## Abstract

### 1. Introduction

Problems requiring choosing an action under uncertainty are rife in all areas of human endeavour. For many problems, actions may be chosen sequentially, allowing the agent to learn from the outcome of early choices to improve later ones.

A widely used framework for sequential decision making is the multi-armed bandit. In the classic multi-armed bandit setting there is a finite set of available actions, each associated with a distribution over rewards which is unknown but stationary. At each timestep the agent selects an action and receives a reward sampled i.i.d from the corresponding reward distribution. The performance of bandit algorithms is described by the regret: the difference in the expected reward obtained by the algorithm and the reward that could be obtained if the optimal action was selected at every timestep.

An alternate approach to selecting actions is causal inference. Frameworks for causal inference provide a mechanism to specify assumptions that allow observational distributions over variables to be mapped to interventional ones. This allows an agent to predict the outcome of an action based on non-experimental data. This approach is common in social science, demography, and economics where explicit experimentation may be difficult. For example, predicting the effect of changes to childcare subsidies on workforce participation or school choice on student grades.

We take a first step towards unifying these approaches by considering a variant of the stochastic multi-armed bandit problem where we have prior knowledge of the causal structure governing the available actions.

A natural way to connect the causal framework with the bandit setting is to model the problem as a causal directed acyclic graph. Each possible assignment of variables to values is an action (bandit arm). The reward could be a general function of the action selected and the final state of the

graph. However for simplicity, we will consider the reward to be the value of a single specified node minus the cost of the selected action. The number of actions grows exponentially with the number of variables in the graph, making it important to use algorithms that take account of the graph structure to reduce the search space.

Problems framed in this way take on characteristics of different bandit settings depending on the assumptions we make about what subset of actions can be taken, what variables are observable and whether they are observed before or after an action is selected. If feedback is received only on the reward node then the do-calculus can be applied to eliminate some actions immediately, before any experiments are performed and then a standard bandit algorithm can be run on the remaining actions.

If we receive feedback on additional nodes the problem can be more interesting. In addition to being able to eliminate some actions prior to sampling any data as in the previous case, taking one action may give us some information on actions that were not selected.

We consider a bandit problem where the actions and reward are represented by a specific causal graph that demonstrates this interesting structure. We develop an algorithm to leverage the information provided by this structure and demonstrate it substantially outperforms standard bandit algorithms applied to the same problem where the number of actions is large.

There has been substantial recent work into extending bandit algorithms to incorporate additional assumptions and deal with more complex feedback structures. Algorithms with strong guarantees have been developed for linear bandits [], generalized linear bandits, gaussian process bandits [], etc. There is also an active line of research into bandits with feedback defined by a graph. Actions are modelled as nodes in the graph and the agent observes rewards for each action connected to the selected action []. The novelty of our work is that we assume prior knowledge of the causal structure but not the functional form of the relationship between variables.

Partial monitoring is a very general framework for for decoupling the feedback from the action and reward. It can be used to classify problems into one of four categories, trivial with no regret, easy with  $R_T = \tilde{\Theta}(\sqrt{T})$ , hard

with  $R_T = \Theta(T^{2/3})$  and hopeless with  $R_T = \Omega(T)$  (?). Partial monitoring algorithms yield results that are optimal with respect to the horizon  $T$  but not other parameters, such as  $K$ , which is the key focus of incorporating causal structure.

ALSO NEED TO MENTION ANY OTHER COMBINATIONS OF BANDITS+CAUSAL (eg the Elias NIPS paper and Generalized Thompson Sampling paper)

Key to Elias' paper is: observing the action an agent would take if it were allowed to make its natural choice can provide some information about hidden confounders that influence both the reward and the choice of action. Therefore, incorporating an agent's natural choice as context may outperform a standard bandit that does not use that context. (Note: even in the presence of hidden confounders, including the agent's natural choice as context only may improve the results. It is easy to come up with a counter example in which it does not).

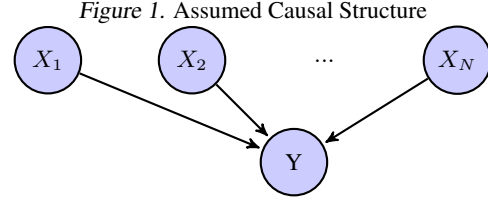
## 2. Problem setup

Assume we have a known causal model with binary variables  $\mathbf{X} = \{X_1 \dots X_N\}$  that independently cause a target variable of interest  $Y$ , figure 1. We can run sequential experiments on the system, where at each timestep  $t$  we can either do nothing,  $do()$ <sup>1</sup>, or select a single variable,  $X_i$ , on which to intervene,  $do(X_{t,i} = J_t)$ , and subsequently observe the complete result,  $(\mathbf{X}_t, Y_t)$ .

As an example, consider a farmer wishing to optimize the yield of her crop. She can invest in a green house to control temperature, a watering system to control soil moisture, fertilizers to set soil nutrients, etc. We assume only a single intervention is feasible due to cost and that each of these variables are independent of one-another (this may not always be the case - temperature could be related to rainfall for example). After having selected which variable to control, she plants her crops and observes the values of the remaining input variables and the yield. This repeats across many growing seasons, and the goal is to maximize the total cumulative yield.

Let  $\mathbf{q} \in [0, 1]^N$  be a fixed vector where  $q_i = P(X_i = 1)$ . In each time-step  $t$  up to a known end point  $T$ :

<sup>1</sup>A note on notation: In the bandit community it is implicit that algorithms selecting actions are intervening in the system. So it is sufficient to index actions according to the variable and value. However, in causal graphs, it is essential to differentiate observing (or conditioning) on a variable taking a certain value, from intervening to set that variable. Although in the specific causal graph we consider, observation and intervention are the same, we deliberately introduce the do-notation (?) that makes this distinction clear so as to help provide a bridge between the bandit and causal inference communities.



1. The learner chooses  $I_t \in \{0 \dots N\}$  and  $J_t \in \{0, 1\}$ , corresponding to setting the variable  $X_{I_t} = J_t$ , also denoted  $do(X_{I_t} = J_t)$ . Selecting  $I_t = 0$  represents not intervening on any variable and simply observing, also denoted  $do()$ .
2. The remaining variables take their values stochastically according to their natural probabilities  $\mathbf{q}$

$$X_{i,t} = \begin{cases} \sim \text{Bernoulli}(q_i) & \text{if } i \neq I_t \\ J_t & \text{otherwise.} \end{cases}$$

3. The learner receives reward  $Y_t \sim \text{Bernoulli}(r(\mathbf{X}_t))$  where  $r : \{0, 1\}^N \rightarrow [0, 1]$  is unknown and arbitrary.

The expected reward of taking action  $i, j$  is  $\mu_{i,j} = \mathbb{E}[r(\mathbf{X}) | do(X_i = j)]$ . The optimal reward and action are denoted  $\mu^*$  and  $(i^*, j^*)$  respectively, where  $(i^*, j^*) = \arg \max_{i,j} \mu_{i,j}$  and  $\mu^* = \mu(i^*, j^*)$ . The  $n$ -step cumulative expected regret is

$$R_n = \mathbb{E} \left[ \sum_{t=1}^n (\mu^* - \mu_{I_t, J_t}) \right].$$

The problem could be treated as a classical multi-armed bandit with  $K = 2N + 1$  arms, yielding a regret  $\mathcal{O}(\sqrt{TN})$ . However, this doesn't leverage the side information induced by the causal structure.

## 3. Causal Bandits

Given the assumed causal structure, the probability of  $Y_t$  given we intervene to set  $X_{i,t} = j$  is the same as if we don't set any variables and observe  $X_{i,t} = j$ .

$$P(Y_t | do(X_{i,t} = j)) = P(Y_t | do(), X_{i,t} = j) \quad (1)$$

This follows from application of the do-calculus (?) to our specific causal graph. It is not the case in general. For example if there was a variable  $X'$  that caused both  $X_i$  and  $Y$  (or  $X_i$  and any other variable  $X_l$ ), that would introduce

a backdoor path from  $X_i \rightarrow Y$  and we would have to condition on  $X'$  to derive the interventional distribution of  $Y$  from the observational one.

Probably should define causal model and back-door rule to properly show this.

We can also learn about the reward for intervening on one variable from rounds in which we actually set a different variable.

$$P(Y_t | do(X_{i,t} = j)) = \sum_{j'} P(Y_t | do(X_{l,t} = j'), X_{i,t} = j) P(X_{l,t} = j') \quad (2)$$

We propose a simple explore-exploit based algorithm that leverages (1). Without loss of generality, we assume  $q_i \in [0, \frac{1}{2}]$  and  $q_1 \leq q_2 \leq \dots \leq q_N$ .

---

**Algorithm 1** Causal Explore-Exploit

---

**Input:**  $T, \mathbf{q}$

Let  $m = \min \{1 \leq i \leq N : q_{i+1} \geq \frac{1}{i}\}$

Let  $h = T^{2/3} m^{1/3} \log(TK)^{1/3}$

Let  $A = \{(i, j) : i \leq m, j = 1\}$  be the set of infrequently observed arms

**for**  $t = 1$  **to**  $h/2$  **do**

    Choose the action  $do()$  and observe  $\mathbf{X}_t$  and  $r_t$

**end for**

Compute for all arms  $(i, j) \notin A$ :

$$\hat{\mu}_{i,j} = \frac{2}{h} \frac{\sum_{t=1}^{h/2} \mathbb{1}\{X_{i,t} = j\} r_t}{q_i^j (1 - q_i)^{1-j}}$$

**for**  $(i, j) \in A$  **do**

**for**  $t' = 1$  **to**  $h/2m$  **do**

        Choose the action  $do(X_{i,t'} = j)$  and observe  $r_t$

**end for**

    Compute  $\hat{\mu}_{i,j} = \frac{2m}{h} \sum_{t'=1}^{h/2m} \mathbb{1}\{X_{i,t'} = j\} r_{t'}$

**end for**

Compute  $(\hat{i}^*, \hat{j}^*) = \arg \max_{(i,j)} \hat{\mu}_{i,j}$

**for**  $t = h$  **to**  $T$  **do**

    Choose the action  $do(X_{\hat{i}^*,t} = \hat{j}^*)$

**end for**

---

**Theorem 1.** Define  $m = \min \{1 \leq i \leq N : q_{i+1} \geq \frac{1}{i}\}$ . Then algorithm 1 satisfies

$$R(T) \in \mathcal{O} \left( T^{2/3} m^{1/3} \log(KT)^{1/3} \right).$$

The lower bound for the standard bandit problem is  $R_t \in \Omega(\sqrt{TK})$  (?). Comparing these results shows exploiting the extra information provided by the causal structure

should outperform standard bandit algorithms when the number of arms is large,  $K \gg m^{2/3} T^{1/3}$ . The parameter  $m$  summarizes the vector  $\mathbf{q}$ , and represents the number of actions that occur rarely naturally and thus must be explicitly explored. If  $q_1, \dots, q_N = 0$ , the problem is completely unbalanced and  $m = N$ . If  $q_1, \dots, q_N = \frac{1}{2}$ , the problem is completely balanced and  $m = 1$ .

Algorithm 1 tries to learn the rewards for all the arms during an exploration phase  $h$  and then picks the arm with the highest empirical mean for all remaining timesteps. During its exploration phase, it learns all the frequently occurring actions by observation and the remaining, infrequently occurring actions, by explicitly playing them. This leads to Chernoff type high probability bounds on the difference between the empirical and true rewards for all arms of the form  $P(\hat{\mu}_{i,j} - \mu_{i,j} > D) \leq e^{-hD^2/m}$ . By choosing optimal values for  $D$  and  $h$  we obtain the regret bound given in theorem 1. A full proof is given in the supplementary materials.

Algorithm 1 relies on  $\mathbf{q}$  and thus  $m$  being known. We now consider the case where  $\mathbf{q}$  is unknown. We begin by considering the simple regret, defined as the expected difference between the mean payoff of the optimal action and that of the action estimated to be optimal in  $T$ th round.

$$R^{\text{simple}}(h) = \mathbb{E} [\mu^* - \max_{(i,j)} \hat{\mu}_{i,j}(h)] \quad (3)$$

**Theorem 2.** Define  $m = \min \{1 \leq i \leq N : q_{i+1} \geq \frac{1}{i}\}$ . Then algorithm 2 satisfies

$$R^{\text{simple}}(T) \in \mathcal{O} \left( \sqrt{\frac{m}{T} \log \left( \frac{NT}{m} \right)} \right).$$

The simple regret for a  $K$ -armed bandit is lower bounded by  $\mathcal{O}(\sqrt{K/T})$  (?). By utilizing the causal structure the dependence on  $K$  is reduced to a dependence on  $m$ , which may be much smaller. The intuition behind algorithm 2 is that we can use data collected during rounds where we select the  $do()$  action to estimate  $m$  and determine which arms we must play explicitly.

For the bandit problem we need to collect some data to estimate  $m$  before determining if we should use the causal explore-exploit algorithm or a standard bandit algorithm.

**Theorem 3.** Define  $m = \min \{1 \leq i \leq N : q_{i+1} \geq \frac{1}{i}\}$ . Then algorithm 3 satisfies

$$R(T) \in \mathcal{O} \left( T^{2/3} m^{1/3} \log(KT)^{1/3} \right).$$

**Algorithm 2** Causal Best Arm Identification

---

```

1: Input:  $T, N$ 
2: for  $t \in 1, \dots, T/2$  do
3:   Choose the action  $do()$  and observe  $\mathbf{X}_t$  and  $r_t$ 
4: end for
5: Compute for all  $i \in \{1, \dots, N\}$  and  $j \in \{0, 1\}$ :

$$\hat{\mu}_{i,j} = \frac{\sum_{t=1}^{T/2} \mathbb{1}\{X_{i,t} = j\} r_t}{\sum_{t=1}^{T/2} \mathbb{1}\{X_{i,t} = j\}}.$$

6: Compute  $\hat{q}_i = \frac{2}{T} \sum_{t=1}^{T/2} X_{i,t}$ 
7: Compute  $\hat{s}_i = \min\{\hat{q}_i, 1 - \hat{q}_i\}$ 
8: Compute  $\hat{s}' = \text{sorted}(\hat{s}) : \hat{s}'_1 \leq \hat{s}'_2 \leq \dots \leq \hat{s}'_N$ 
9: Compute  $\hat{m} = \min\{1 \leq i \leq N : \hat{s}'_{i+1} \geq \frac{1}{2}\}$ 
10:  $i'(i)$  = the index of  $\hat{s}_i$  in  $\hat{s}'$ 
11: Compute  $A$  as the subset of infrequently observed arms
 $\{(i, j) : i'(i) \leq \hat{m}, j = \mathbb{1}\{\hat{q}_i \leq \frac{1}{2}\}\}$  with  $|A| = \hat{m}$ 
12: for  $(i, j) \in A$  do
13:   for  $t \in 1, \dots, T/2\hat{m}$  do
14:     Choose action  $do(X_{i,t} = j)$  and observe  $r_t$ 
15:   end for
16:   Recompute  $\hat{\mu}_{i,j} = \frac{2\hat{m}}{T} \sum_{t=1}^{T/2\hat{m}} r_t(X_{i,t} = j)$ 
17: end for

```

---

**3.1. Lower Bounds****4. Experiments****5. Discussion**

In our algorithm, we have only used the side information provided by the  $do()$  action about other actions. Since the  $do()$  action fully reveals the value of alternate actions we could have incorporated this information via the graph feedback model (?), where at each timestep the feedback graph  $G_t$  is selected stochastically, dependent on  $\mathbf{q}$ , and revealed after an action has been chosen. The feedback graph is distinct from the causal graph. A link  $A \rightarrow B$  in  $G_t$  indicates that selecting the action  $A$  reveals the reward for action  $B$ . For this specific problem,  $G_t$  will always be a star graph with the action  $do()$  connected to half the remaining actions. The Exp3-IX algorithm (?) was developed for the adversarial version of this problem and has regret  $\mathcal{O}(\sqrt{\bar{\alpha}T})$ , where  $\bar{\alpha}$  is the average independence number of  $G_t$ . In our case  $\bar{\alpha} = \frac{N}{2}$  so we again obtain the regret of the standard bandit algorithm. The issue here is that a malicious adversary can select the same graph each time, such that the rewards for half the arms are never revealed by the informative action. This is equivalent to a, nominally, stochastic selection of feedback graph where  $\mathbf{q} = \mathbf{0}$

(?) consider a stochastic version of the graph feedback problem, but with a fixed graph available to the algorithm

Figure 2. Final regret versus number of variables  $N$  for UCB with  $\alpha = 2$ , Causal-Explore-Exploit with  $m = 2$  and with  $m = N$  and horizon  $T = 10,000$ . Error bars show standard deviation over 100 simulations. The regret for UCB grows linearly with the number of variables, whilst for Causal-Explore-Exploit with fixed  $m$ , the growth is sub-logarithmic.

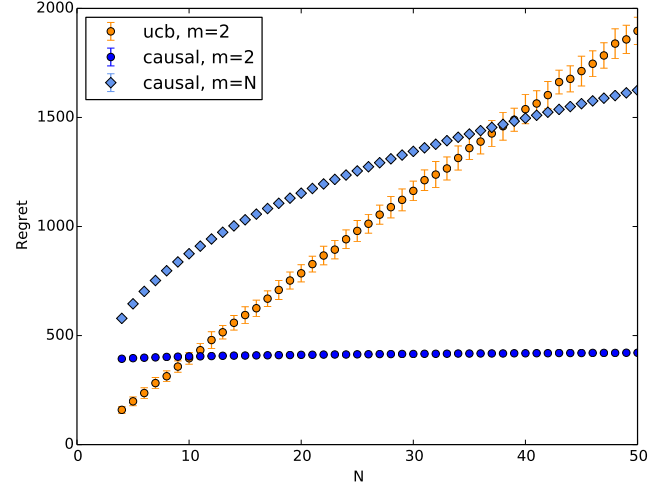
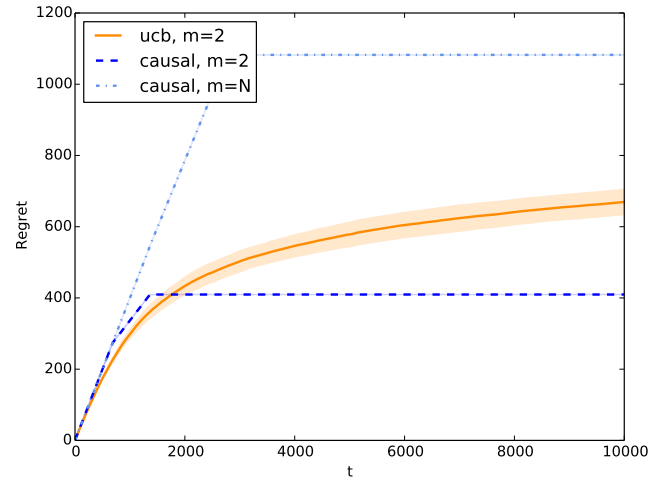


Figure 3. Cumulative regret over time for  $N = 17$  for UCB with  $\alpha = 2$ , Causal-Explore-Exploit with  $m = 2$  and Causal-Explore-Exploit with  $m = N$ . Shaded region shows standard deviation over 100 simulations. The Causal-Explore-Exploit algorithm incurs linear regret during the exploration phase, after which it selects the optimal arm with high probability. For  $m = 2$ , we have  $K \sim m^{2/3}T^{1/3}$  and see that we are in the regime in which Causal-Explore-Exploit outperforms UCB.



**Algorithm 3** Bandit Regret Algorithm

---

```

1: Input:  $T, N$ 
2:  $\delta = \frac{1}{T^{1/3}}$ 
3:  $T_1 = 48N \log(4N/\delta)$ 
4: Run algorithm 2 to line 11 with input  $T_1, N$ .
5: if  $\hat{m} > \frac{N^{3/2}}{\sqrt{T}}$  then
6:   Switch to the standard UCB algorithm.
7: else
8:    $h = T^{2/3} \hat{m}^{1/3} \log(TK)^{1/3}$ 
9:   Run algorithm 2 with input  $h, N$ .
10:  Compute  $(\hat{i}^*, \hat{j}^*) = \arg \max_{(i,j)} \hat{\mu}_{i,j}$ 
11:  for  $t = h$  to  $T$  do
12:    Choose the action  $do(X_{\hat{i}^*, t} = \hat{j}^*)$ 
13:  end for
14: end if

```

---

before it must select an action. In addition, their algorithm is not optimal for all graph structures and fails, in particular, to provide improvements for star like graphs as in our case. (?) improve the dependence of the algorithm on the graph structure but still assume the graph is fixed and available to the algorithm before the action is selected.

More generally, assuming causal structure creates more complex types of side information, such as that shown in equation 2. In this case, selecting one action does not fully reveal an alternate action but provides some information towards an estimate. The quality of the estimate notably depends not only on the number of times that action was selected. For example, to get a good estimate for  $X_1 = 1$  by intervening on  $X_2$  requires us to sample both  $X_2 = 0$  and  $X_2 = 1$ , in proportions dependent on  $q_2$ . This more complex side information does not fit within the graph feedback framework.

## 6. Future Open Questions

- Known but arbitrary structure
- Learning structure then exploiting

## 7. Conclusion