

---

# Bandit Data-driven Optimization: AI for Social Good and Beyond

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The use of machine learning (ML) systems in real-world applications entails  
2 more than just a prediction algorithm. AI for social good (AI4SG) applications,  
3 and many real-world ML tasks in general, feature an iterative process that joins  
4 prediction, optimization, and data acquisition in a loop. We introduce bandit  
5 data-driven optimization, the first iterative prediction-prescription framework to  
6 combine the advantages of online bandit learning and offline predictive analytics. It  
7 offers a flexible setup to reason about unmodeled policy objectives and unforeseen  
8 consequences. We propose PROOF, the first algorithm for this framework and show  
9 that it achieves no-regret. PROOF achieves superior performance over existing  
10 baseline in numerical simulations.

## 11 1 Introduction

12 The success of modern ML often does not directly translate into a perfect solution to a real-world  
13 AI4SG problem [17], or more generally, data-driven policy-making. One obvious reason is that  
14 supervised learning focuses on prediction, yet real-world problems, by and large, need prescription.  
15 For example, rather than predict which households' water pipes are contaminated (labels) using  
16 construction and demographic information (features), municipal officials need to know how to  
17 schedule inspections and replacements (interventions) [2]. The common practice is a two-stage  
18 procedure, as shown in Figure 1a. After training an ML prediction model, the user makes prescriptive  
19 decisions based on some optimization problem or even simple heuristic which takes the prediction  
20 output as parameters. The training objective and the optimization objective are completely separate,  
21 which means it is hard to control the final prescription quality. In an emerging line of work on (one-  
22 shot) data-driven optimization, the learning problem is made aware of the downstream optimization  
23 objective through its loss function, and hence mitigating this issue [7, 11]. We illustrate this in  
24 Figure 1b.

25 However, this is still far from the complete picture of a real-world problem. Figure 1c shows a typical  
26 workflow in many AI4SG applications. After getting data from the collaborating organization, the  
27 researcher trains a predictive model and then, based on it, makes an intervention recommendation.  
28 The workflow does not stop here, though. After the organization implements the recommended  
29 intervention, it collects more data points. Using these additional data, the researcher then updates  
30 the predictive model and makes a new intervention recommendation to be implemented, so on and  
31 so forth, resulting in an iterative process. The guiding principles of the various components in this  
32 process are often not aligned. Without a rigorous, integrated framework to guide the procedure, this  
33 could lead to operation inefficiency, missed expectations, dampened initiatives, and new barriers of  
34 mistrust which are not meant to be.

35 Why is such an iterative process necessary? First, many "social good" domains do not have the  
36 luxury of millions of training examples. A small dataset leads to inaccurate predictions and hence

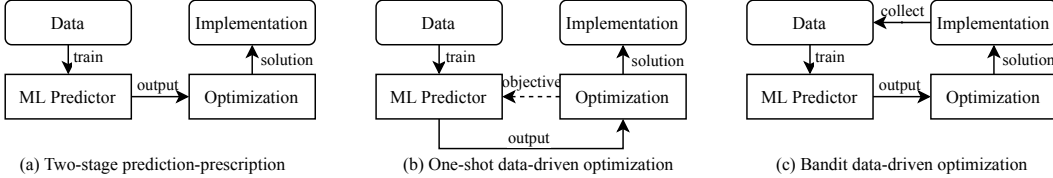


Figure 1: Paradigms of how ML systems are used in realistic policy-making.

37 suboptimal decisions. Second, too often the initial dataset has some default intervention embedded,  
 38 while the project’s goal is to find the optimal intervention. It is necessary to at least try out some  
 39 interventions and collect data under them. Third, the communication gap between researchers and  
 40 practitioners makes it difficult to formulate the right objective at the beginning. Fourth, interventions  
 41 may have unexpected consequences, thus the inherent impossibility of fully modeling the problem in  
 42 one shot.

43 We propose the first iterative prediction-prescription framework, which we term as bandit data-driven  
 44 optimization. Bandit data-driven optimization combines the relative advantages of both online bandit  
 45 learning and offline predictive analytics. We achieve this with our algorithm PRedict-then-Optimize  
 46 with Optimism in Face of uncertainty (PROOF). PROOF is a modular algorithm which can work  
 47 with a variety of predictive models and optimization problems. Under specific settings, we formally  
 48 analyze its performance and show that PROOF achieves no-regret. In addition, we propose a variant  
 49 of PROOF which can handle the scenario where the intervention affects the data distribution, which  
 50 also enjoys no-regret. Finally, we use numerical simulations to show that PROOF achieves superior  
 51 performance than a pure bandit baseline.

52 We emphasized AI4SG as the motivation and application domain of bandit data-driven optimization,  
 53 because we have experienced these problems first-hand in our applied work on AI4SG. However, none  
 54 of the reasons above is exclusive to AI4SG. Many data science projects in the real world, whether  
 55 explicitly for social good or not, fit these characterizations and can benefit from such a framework. In  
 56 this sense, our work can be viewed as a general framework for data-driven policy-making.

## 57 2 Related Work

58 There are many paradigms in which ML is used in a real-world problem, yet the paradigm illustrated  
 59 in Figure 1c is certainly a common one. To our knowledge, there is surprisingly no existing work  
 60 that rigorously studies this procedure. We propose bandit data-driven optimization as the first step  
 61 towards formalizing and improving the way machine learning is used in AI4SG projects. That said,  
 62 several research topics also concern themselves with similar problems in the real-life usage of ML  
 63 systems. We introduce them below and compare each of them with bandit data-driven optimization.

64 First, (one-shot) data-driven optimization is an emerging line of research in the operations research  
 65 literature. Given a dataset consisting of features  $x_1, \dots, x_n$  and labels  $c_1, \dots, c_n$ , the task is to find the  
 66 action  $w^*$  that maximizes the expected utility given some feature  $x$ , i.e.  $w^* = \arg \max_w \mathbb{E}_{c|x} [p(c, w)]$ .  
 67 There are two major approaches. The first one comes from the stochastic programming perspective [7,  
 68 6]. For example, Bertsimas and Kallus fuse the prediction and optimization by transforming the  
 69 optimization objective into a weighted combination of cost, where the weight, determined by some  
 70 ML algorithm, represents the similarity between the target feature  $x$  and each of the feature  $x_i$ ’s  
 71 in the dataset. Another popular approach is generally referred to as the predict-then-optimize  
 72 framework [13, 11, 14]. The idea is to fit an ML predictor  $f$  from the feature  $x$  to label  $c$ , and then  
 73 use the prediction  $c = f(x)$  in the optimization problem. To connect the training and optimization,  
 74 typically the loss function in the ML problem is modified to reflect the downstream optimization  
 75 objective. Compared to bandit data-driven optimization, this entire literature assumes that the  
 76 optimization objective is known a priori and does not consider multi-period settings. This limits its  
 77 applicability in reality.

78 Speaking of multi-period settings, contextual bandit is a well-studied online decision-making model  
 79 which is very relevant to our framework [5, 15]. At each time step  $t$  of the contextual bandit,  
 80 we receive feature  $x_t$ , pick an action  $w_t$ , and receive reward whose expectation  $q(x_t, w_t)$  is some

---

**Procedure 1: BANDIT DATA-DRIVEN OPTIMIZATION**

---

- 1 Receive initial dataset  $\mathcal{D} = \{(x_i^0, c_i^0; w_i^0)_{i=1, \dots, n}\}$  from distribution  $D$  on  $(X, C)$ .
  - 2 **for**  $t = 1, 2, \dots, T$  **do**
  - 3     Using all the available data  $\mathcal{D}$ , train ML prediction model  $f_t : X \rightarrow C$ .
  - 4     Given  $n$  feature samples  $\{x_i^t\} \sim D_x$ , choose interventions  $w^t = \{w_i^t\}$  for each individual  $i$ .
  - 5     Receive  $n$  labels  $\{c_i^t\} \sim D(w_i^t)_{c|x_i^t}$ . Add  $\{(x_i^t, c_i^t; w_i^t)_{i=1, \dots, n}\}$  to the dataset  $\mathcal{D}$ .
  - 6     Get cost  $u_t = u(C^t, w^t) = \sum_i p(C_i^t, w_i^t) + q(w_i^t) + \eta$ , where  $\eta \sim N(0, \sigma^2)$ .
- 

81 unknown function of  $x_t$  and  $w_t$ . In fact, contextual bandit is even more general than bandit data-driven  
82 optimization, because we could simply ignore the label  $c$ , forego the training of a predictive model,  
83 and let the bandit algorithm pick an action. However, by doing so, we effectively give up all the  
84 valuable information in the historical data, finding the optimal action in a purely online fashion.  
85 Although contextual bandit algorithms often achieve no-regret and have been used for high-frequency  
86 decision-making [16], they are often impractical in AI4SG applications. It would hardly be acceptable  
87 to any AI4SG stakeholders that our algorithm only guarantees good results after using it for, say, 10  
88 years, while the algorithm chooses not to use the dataset already available. That being said, bandit  
89 provides a proper setting for sequential decision making under uncertainty and bandit algorithms like  
90 LinUCB [9, 8, 1] play a central role in designing algorithms for bandit data-driven optimization.

91 Also related to our problem is offline policy learning [19, 10, 4], which shares similar goals with  
92 bandit data-driven optimization. It aims at finding the best policy  $\pi(w|x)$ , which is a distribution over  
93 actions for each given feature, using only historical action records. Its main advantage over its online  
94 counterpart, i.e. contextual bandit, is that it does not need even a single online trial, and hence is much  
95 easier to convince the stakeholder to adopt. However, this advantage comes with the assumption that  
96 the historical data has many different actions attempted. This assumption often fails to hold, at least  
97 in the AI4SG projects we have worked on. Furthermore, most of this literature focus on the binary  
98 action setting and it, similar to contextual bandits, does not explicitly use the feature/label dataset.

### 99 3 Bandit Data-driven Optimization

100 We describe the formal setup of bandit data-driven optimization in Procedure 1. On Line 1, we  
101 receive an initial dataset  $\mathcal{D}$  of size  $n$ , consisting of features  $x$  and labels  $c$ . Each feature vector  $x_i^0$  is  
102 drawn i.i.d. from an unknown distribution  $D_x$ . Each label vector  $c_i^0$  is independently drawn from an  
103 unknown conditional distribution  $D(w_i^0)_{c|x_i^0}$ . The conditional distribution is indexed by intervention  
104  $w$ , which means different interventions could potentially lead to different data distributions. In this  
105 case,  $w_i^0$  is the default intervention which has been in-place with the data collector. On Line 3,  
106 we use all the data collected so far to train a machine learning model  $f_t$ , a mapping from features  
107  $X \subseteq \mathbb{R}^m$  to labels  $C \subseteq \mathbb{R}^d$ . On Line 4, we draw a new sample of features  $\{x_i^t\}_{i=1}^n$ . Then, we  
108 select an intervention  $w_i^t \in W$  for each individual  $i$ . We have a known loss function  $p(c, w)$  which  
109 represents our modeling effort and domain knowledge. We want to minimize its expectation given the  
110 features  $x_i^t$ . On Line 5, we commit to intervention  $w^t$  and then get the labels  $\{c_i^t\}_{i=1}^n$  corresponding  
111 to the features drawn earlier. Subsequently, on Line 6, we incur a cost  $u_t$ . The cost consists of our  
112 known optimization objective  $p(\cdot)$  and an unknown cost  $q(\cdot)$  that represents all the unmodeled policy  
113 objectives and the unintended consequences. It comes with a random noise.

114 A few remarks are in order. First, a common question is why do we need an iterative procedure at  
115 all? From the real-world application perspective, data are often hard to acquire in AI4SG domains,  
116 so frequently we have no choice but to collect data as we roll out the intervention programs. Even  
117 when we seem to have enough data, it is still wise to keep collecting data and updating our model,  
118 because the data distribution might change over time, and continuous engagement is an important  
119 factor in AI4SG projects. The iterative procedure is also essential for exploration purpose. Bandit  
120 data-driven optimization contains two types of exploration-exploitation trade-offs. The historical data  
121 often contains only a single default intervention. If we do not explore and implement some other  
122 interventions, we will never learn how good (or bad) they are. The second exploration arises when  
123 dealing with the unknown cost component  $q(\cdot)$ , which is essentially a bandit problem in itself. More  
124 discussion on  $q(\cdot)$  follows below.

125 Second, this form of loss – a known part  $p(\cdot)$  and an unknown part  $q(\cdot)$  – is a realistic compromise of  
 126 two extremes. AI4SG researchers often spend a considerable amount of time communicating with  
 127 domain experts to understand the problem. It would go against this honest effort to eliminate  $p(\cdot)$  and  
 128 model the process as a pure bandit problem. On the other hand, even when we take our completed  
 129 work to the deployment stage, there will still be unmodeled objectives. Thus, it would be too arrogant  
 130 to eliminate  $q(\cdot)$  and to pretend that anything that does not go according to the plan is noise. The  
 131 unknown  $q(\cdot)$  is also our conscious acknowledgement that any intervention recommended by AI4SG  
 132 projects may have unintended consequences. We believe that having both  $p(\cdot)$  and  $q(\cdot)$  is a faithful  
 133 first step towards capturing the nature of AI4SG work.

134 Let us use two AI4SG examples to illustrate how bandit data-driven optimization captures real-world  
 135 ML policy-making workflows.

136 **Food Rescue** A food rescue (FR) organization receives food donations from restaurants and grocery  
 137 stores and connects them to low-resource community organizations. Dispatchers at FR would post  
 138 the donor and recipient information on their mobile app, and some volunteer would claim the rescue  
 139 and pick up and deliver the donations. To ensure that each rescue gets claimed, the dispatcher sets  
 140 mobile app push notification time for each rescue, which is the intervention  $w$ , to alert the volunteers  
 141 of available rescues. This decision is dependent on how likely a rescue will be claimed. Thus, we can  
 142 develop a machine learning model which uses features  $x$  of a rescue, e.g. donor/recipient location,  
 143 weather, and time of day, to predict the probability that a rescue will be claimed in 10 minutes, 30  
 144 minutes, etc. (label  $c$ ). The optimization objective  $p(c, w)$  is to minimize the expected wait time,  
 145 while still guaranteeing not a lot of push notifications are sent. After we select a  $w$  for a rescue, we  
 146 observe the label and this data point will be used for training before the next rescue trip comes along.  
 147 The cost to the FR may include factors other than the rescue claim rate, e.g. gain/loss of registered  
 148 volunteers as a result of push notification scheme. Of course we know this now, but the  $q(\cdot)$  cost  
 149 could capture similar factors which we do not know yet. For more background about the food rescue  
 150 operation, the reader may refer to the work of Shi et al. [18].

151 **Anti-poaching** Wildlife poaching is a pressing problem in many parts of the world. Given that  
 152 poachers are often experienced and shrewd, wildlife rangers need to patrol the vast area of a wildlife  
 153 park intelligently. Indeed, this battle between rangers and poachers have received considerable interest  
 154 from the computer science community. This is typically studied as a Stackelberg security game,  
 155 which is essentially an optimization problem to find the best ranger patrolling strategy (intervention  
 156  $w$ ). To solve the optimization problem, we need to know the poacher’s behavior pattern, for which  
 157 we will develop an ML model. The ML model takes as input the location, geographic feature, animal  
 158 density, etc., of a small patch of land (feature  $x$ ) and predicts the likelihood that the poacher will poach  
 159 at this location (label  $c$ ). Given this likelihood, we can select the optimal ranger patrolling strategy  
 160  $w$  which minimizes the likelihood of successful poaching. Every time a poacher or their belonging  
 161 is caught, we can add a new data point to the dataset. The cost to the rangers may include factors  
 162 other than the poached animals themselves, such as the crime rate in the region, as the enforcement of  
 163 anti-poaching might lead to the increase of other crimes.  $q(\cdot)$  could capture similar factors which we  
 164 do not know yet. For more background about the game-theoretic work on anti-poaching, the reader  
 165 may refer to the work of Fang et al. [12].

## 166 4 PRedict-then-Optimize with Optimism in Face of uncertainty (PROOF)

167 We propose the first algorithm for the bandit data-driven optimization, PRedict-then-Optimize with  
 168 Optimism in Face of uncertainty (PROOF), which is shown in Algorithm 2.

169 Unless otherwise specified, we work with the following setting. The data points are drawn from  
 170  $X \times C$  where  $X \subseteq \mathbb{R}^m$  and  $C \subseteq \mathbb{R}^d$ . We assume all  $x \in X$  has  $l^2$ -norm bounded by constant  $K_X$ ,  
 171 and the label space  $C$  has  $l^1$ -diameter  $K_C$ . The action space  $W$  could be either discrete or continuous,  
 172 but is bounded inside the unit  $l^2$ -ball in  $\mathbb{R}^d$ . We specify the data distribution by an arbitrary marginal  
 173 distribution  $D_x$  on  $X$  and a conditional distribution such that  $c = f(x) + \epsilon$  where  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ ,  
 174 for some function  $f$ . Thus,  $\mathbb{E}_{c|x}[c] = f(x)$ . Of course,  $f$  is unknown to us and needs to be learned.  
 175 To begin with, we assume  $f \in \mathcal{F}$  comes from the class of all linear functions with  $f(x) = Fx$ , and  
 176 we use ordinary least squares regression as the learning algorithm. We will relax this assumption  
 177 towards the end of Section 4. The known cost  $p(c, w) = c^\dagger w$  is the inner product of label  $c$  and

---

**Algorithm 2: PROOF: PREDICT-THEN-OPTIMIZE WITH OPTIMISM IN FACE OF UNCERTAINTY**


---

```

1 Initialize:
2   Find a barycentric spanner  $b_1, \dots, b_d$  for  $W$ 
3   Set  $A_i^1 = \sum_{j=1}^d b_j b_j^\dagger$  and  $\hat{\mu}_i^1 = 0$  for all  $i$ 
4 Receive initial dataset  $\mathcal{D} = \{(x_i^0, c_i^0; w_i^0)_{i=1, \dots, n}\}$ .
5 for  $t = 1, 2, \dots, T$  do
6   Using all the available data  $\mathcal{D}$ , train ML prediction model  $f_t : X \rightarrow C$ .
7   Given  $n$  feature samples  $\{x_i^t\} \sim D_x$ , get predictions  $\hat{c}_i^t = f_t(x_i^t)$ .
8   Set confidence ball radius  $\beta^t = \max \left( 128d \log t \log(nt^2/\gamma), \left(\frac{8}{3} \log \left(\frac{nt^2}{\gamma}\right)\right)^2 \right)$ 
9   for  $i = 1, 2, \dots, n$  do
10    Set CB  $B_i^t = \{\nu : \|\nu - \hat{\mu}_i^t\|_{2, A_i^t} \leq \sqrt{\beta^t}\}$ .
11    Choose intervention  $w_i^t$  where  $w_i^t = \arg \min_{w \in W} \min_{\nu \in B_i^t} (\hat{c}_i^t + \nu)^\dagger w$ .
12    Get label  $c_i^t \sim D_{c|x_i^t}$ . Add  $(x_i^t, c_i^t; w_i^t)$  to  $\mathcal{D}$ .
13    Get cost  $u_i^t = (c_i^t)^\dagger w_i^t + \mu^\dagger w_i^t + \eta_i$ , where  $\eta_i \sim N(0, \sigma^2)$ . In particular, let  $u_{oi}^t$  be the
        first term and let  $u_{bi}^t$  be the sum of the second and third term.
14    Update  $A_i^{t+1} = A_i^t + w_i^t (w_i^t)^\dagger$ 
15    Update  $\hat{\mu}_i^{t+1} = (A_i^{t+1})^{-1} \sum_{\tau=1}^t u_{bi}^\tau w_i^\tau$ 

```

---

178 action  $w$ .<sup>1</sup> The unknown cost is  $q(w) = \mu^\dagger w$ , where  $\mu$  is an unknown but fixed vector. Furthermore,  
179 for exposition purpose we will start by assuming that the intervention  $w$  does not affect the data  
180 distribution. We will later remove this assumption and present the algorithm for the general case.

181 PROOF is a delicate integration of the celebrated Optimism in Face of Uncertainty (OFU) frame-  
182 work [9] and the predict-then-optimize framework. It is clear that the unknown cost component  
183  $q(\cdot) + \eta$  forms a linear bandit. For this bandit component, we run an OFU algorithm for each  
184 individual  $i$  with the same unknown loss vector  $\mu$ . The OFU component for each individual  $i$   
185 maintains a confidence ball ( $B_i^t$ ) which is independent of the predict-optimize framework. The  
186 predict-then-optimize framework produces an estimated optimization objective (represented by  $\hat{c}^t$ )  
187 that has nothing to do with the OFU. The two components are integrated together on Line 11 of  
188 Algorithm 2, where we compute the intervention for the current round taking into consideration the  
189 essence of both frameworks.

190 Below, we justify why this algorithm achieves no-regret. First, we state a theorem by Dani et al. [9],  
191 which states that the confidence ball captures the true loss vector  $\mu$  with high probability. Our only  
192 modification is a trivial union bound so that the result holds for all the  $n$  bandits simultaneously.

193 **Lemma 1** (Adapted from Theorem 5 by Dani et al. [9]). *Let  $\gamma > 0$ ,*

$$\mathbb{P}(\forall t, \forall i, \mu \in B_i^t) \geq 1 - \gamma.$$

194 This lemma was proved for the original OFU algorithm. However, the lemma itself does not depend  
195 on the way we choose  $w^t$  at each time step, as long as we computed  $\hat{\mu}$  and  $A$  in the way we did (same  
196 as OFU) on Lines 14-15. That is, it does not matter that we have an additional  $\hat{c}_i^t$  in the optimization  
197 on Line 11, compared to the original OFU. We state another useful result below.<sup>2</sup>

198 **Lemma 2.** *Suppose we use the ordinary least squares regression as the ML algorithm. The prediction*  
199 *error is*

$$\mathbb{E}_{X, \epsilon} \left[ \left\| \mathbb{E}_{c_i^t | x_i^t} [c_i^t] - \hat{c}_i^t \right\|_2 \right] = O \left( \sqrt{\frac{dm}{nt}} \right).$$

200

201 **Theorem 3.** *Assuming we use ordinary least squares regression as the ML algorithm, PROOF has*  
202 *regret  $\tilde{O} \left( n\sqrt{dmT} \right)$  with probability  $1 - \delta$ .*

<sup>1</sup>To avoid confusion, in this paper we use superscript  $\dagger$  to denote matrix and vector transpose.

<sup>2</sup>Please refer to the full paper for complete proofs and additional technical results. [https://www.dropbox.com/s/92nk7z0lycva90h/BD0\\_anon.pdf?dl=0](https://www.dropbox.com/s/92nk7z0lycva90h/BD0_anon.pdf?dl=0)

203 *Proof.* Let  $w_{i*}^t = \arg \min_w (\mathbb{E}_{c_i^t|x_i^t}[c_i^t] + \mu)^\dagger w$ .  $w_{i*}^t$  is the optimal action for individual  $i$  at time  $t$ ,  
 204 and is the benchmark in our regret computation.

205 Fix  $i$ , fix  $t$ . Let  $\tilde{\nu} = \arg \min_{\nu \in B_i^t} (\hat{c}_i^t + \nu)^\dagger w_{i*}^t$ . Because of Line 11, we have

$$(\hat{c}_i^t + \tilde{\nu})^\dagger w_{i*}^t = \min_{\nu \in B_i^t, w \in W} (\hat{c}_i^t + \nu)^\dagger w \leq (\mathbb{E}_{c_i^t|x_i^t}[c_i^t] + \mu)^\dagger w_{i*}^t + (\hat{c}_i^t)^\dagger w_{i*}^t - \mathbb{E}_{c_i^t|x_i^t}[c_i^t]^\dagger w_{i*}^t.$$

206 The inequality above used the fact that  $\mu \in B_i^t$ , by Lemma 1. Thus, we get the per-round regret

$$\begin{aligned} (\mathbb{E}_{c_i^t|x_i^t}[c_i^t] + \mu)^\dagger (w_i^t - w_{i*}^t) &\leq (\mathbb{E}_{c_i^t|x_i^t}[c_i^t] + \mu)^\dagger w_i^t - (\hat{c}_i^t + \tilde{\nu})^\dagger w_i^t + (\hat{c}_i^t)^\dagger w_{i*}^t - \mathbb{E}_{c_i^t|x_i^t}[c_i^t]^\dagger w_{i*}^t \\ &= (\mathbb{E}_{c_i^t|x_i^t}[c_i^t] - \hat{c}_i^t)^\dagger (w_i^t - w_{i*}^t) + (\mu - \tilde{\nu})^\dagger w_{i*}^t \end{aligned}$$

207 We can view the second term is the per-round regret for the bandit part. By Theorem 6 in Dani et  
 208 al. [9], we have

$$\sum_{t=1}^T ((\mu - \tilde{\nu})^\dagger w_{i*}^t)^2 \leq 8m\beta^T \log T$$

209 Using the Cauchy-Schwarz, we get

$$\sum_{t=1}^T (\mu - \tilde{\nu})^\dagger w_{i*}^t \leq \sqrt{8mT\beta^T \log T}$$

210 Thus, the regret of PROOF is

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^n (\mathbb{E}_{c_i^t|x_i^t}[c_i^t] + \mu)^\dagger (w_i^t - w_{i*}^t) \right] &\leq \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^n (\mathbb{E}_{c_i^t|x_i^t}[c_i^t] - \hat{c}_i^t)^\dagger (w_i^t - w_{i*}^t) \right] + n\sqrt{8mT\beta^T \log T} \\ &= O \left( \sum_{t=1}^T \sum_{i=1}^n \mathbb{E} \left[ \left\| \mathbb{E}_{c_i^t|x_i^t}[c_i^t] - \hat{c}_i^t \right\|_2 \right] + n\sqrt{8mT\beta^T \log T} \right) \end{aligned}$$

211 The last step above used Cauchy-Schwartz and the bounded action space assumption. Using Lemma 2,  
 212 we can conclude the total regret is

$$R_T = O \left( \sum_{t=1}^T \sum_{i=1}^n \sqrt{\frac{dm}{nt}} + n\sqrt{8mT\beta^T \log T} \right) = O \left( n\sqrt{dmT} \right)$$

213 The last step (bounding  $\sum_{t=1}^T t^{-1/2}$ ) is by an upper bound on the generalized harmonic numbers,  
 214 which can be found in Theorem 3.2 (b) in the text by Apostol [3].  $\square$

215 We now consider the more general setting where the intervention could affect the label distribution.  
 216 We make the assumption that there are finitely many possible actions. In the full paper, we show that  
 217 an adaptation of PROOF achieves no-regret, and we also consider continuous action space.

218 **Theorem 4.** *Assuming an OLS regression ML problem, an adaptation of PROOF has regret*  
 219  $\tilde{O} \left( n(d|W|)^{1/3} m^{1/2} T^{2/3} \right)$  *with probability  $1 - \delta$ .*

## 220 5 Numerical Simulations

221 In this section, we implement the PROOF algorithm described in Section 4 and show its performance  
 222 on a simulated dataset. Recall that we train an ML predictor  $\hat{f} : X \rightarrow C$  where  $X \subseteq \mathbb{R}^m$  and  
 223  $C \subseteq \mathbb{R}^d$ . For a data point  $(x, c)$ , we want to solve the optimization problem  $w^* = \min_{w \in W} (c + \mu)^\dagger w$ ,  
 224 had we known  $c$  and  $\mu$ , where  $\mu$  is the unknown bandit parameter. In Section 4, we showed that our  
 225 PROOF algorithm can achieve no-regret when the true ML predictor is a linear mapping.

226 We start our numerical experiments with a small-scale setting. We take feature dimension  $m = 20$   
 227 and label dimension  $d = 5$ . At every round we get  $n = 20$  data points. Following the tradition in the  
 228 bandit literature, we assume the bandit reward is bounded in  $[-1, 1]$  and assume the feasible region  
 229  $W$  is the unit  $l_2$ -ball, as it's not the absolute magnitude, but the relative magnitude between the bandit  
 230 and the optimization rewards, that matters. For the true linear map  $F$ , i.e.  $c = Fx + \epsilon$ , we upper

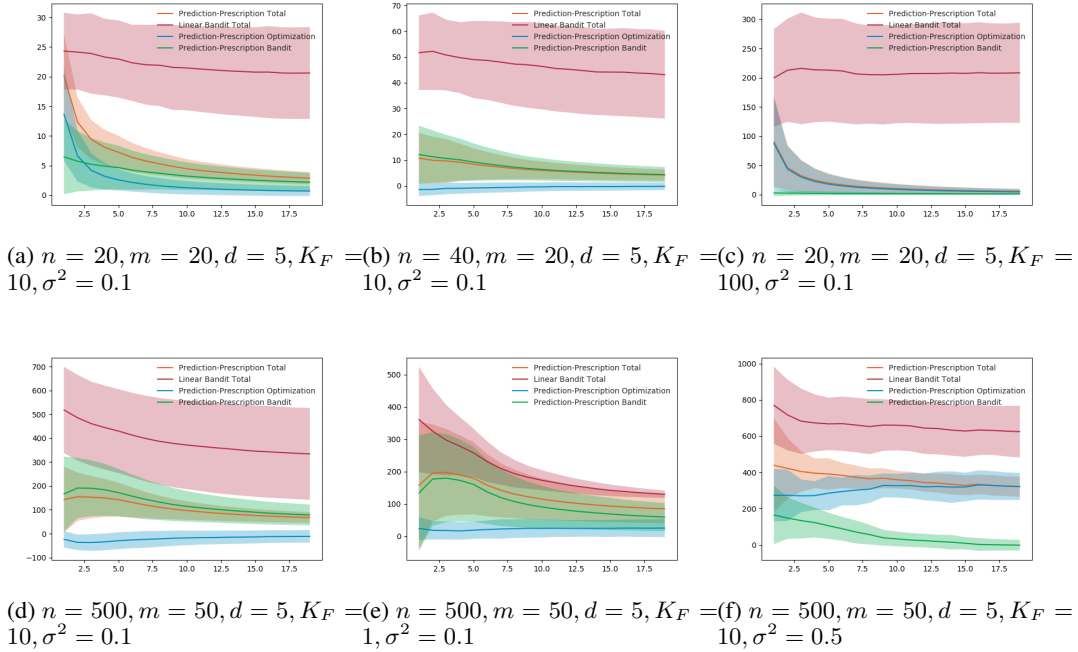


Figure 2: Numerical simulation results of PROOF compared against vanilla linear bandit. All results are averaged over 10 runs with shaded areas representing the standard deviation.

231 bound its  $l_1$  matrix norm at 10. We sample the noise  $\epsilon \sim \mathcal{N}(0, \sigma^2 I_d)$  from a normal distribution  
 232 where  $\sigma^2 = 0.1$ . For the bandit noise, we take  $\eta \sim N(0, 10^{-4})$ . We use ordinary least squares  
 233 regression on all the data points collected so far for the learning part at each iteration of bandit  
 234 data-driven optimization. The optimization at each iteration of PROOF (Line 11 in Algorithm 2) is  
 235 a non-convex bilinear optimization problem. We use the non-convex solver IPOPT to compute a  
 236 heuristic solution to this problem. To compute the regret we also need to find the best action given  
 237 the true reward parameters. This is a convex problem with linear objective and convex quadratic  
 238 constraints, and thus we use Gurobi to solve it. All experiment results are the average over 10 runs.

239 In Alg. 2, there is an explicit expression for the confidence ball’s radius  $\beta^t$ . That is for establishing  
 240 the theoretical regret bound. However, in experiments, we found that, following that formula, the  
 241 radius would be too large ( $\sim 10^4$ ). This makes the algorithm unable to select meaningful action in  
 242 the early rounds, even when the algorithm’s reward estimate  $\hat{\mu}$  is already quite accurate. In practice,  
 243 it is common to select a small value of  $\beta$ , so that the algorithm can quickly concentrate on the correct  
 244 region of interest. We thus set  $\beta^t = 1$ .

245 Under this setting, the problem in theory might be solved simply as a linear bandit problem. The  
 246 expected cost for a fixed action  $w$  is

$$\mathbb{E}_{c,\eta}[(c + \mu)^\dagger w + \eta] = \mathbb{E}[x^\dagger F^\dagger w] + \mu^\dagger w = \mathbb{E}[x^\dagger] F^\dagger w + \mu^\dagger w = \mu^\dagger w,$$

247 because when we generated  $x$ , the distribution has zero mean. This fits in the setting of [9] and thus  
 248 we could feed the total cost in bandit data-driven optimization to their OFU algorithm. Although  
 249 vanilla OFU completely ignores the predict-then-optimize procedure, its regret bound is still the same  
 250 as our PROOF in terms of the order of  $T$ . This brings back the question that we have been repeating  
 251 since the beginning of this paper: if linear (contextual) bandit is a more general framework, why  
 252 should we care about predict-then-optimize at all?

253 We have answered this question with the nature of AI4SG projects and real-world applications of  
 254 ML systems. Here, we can also answer this question using numerical experiments. We show the  
 255 average regret of PROOF as the orange curve in Figure 2, and that of OFU in red. We can decompose  
 256 the average regret of PROOF into the regret of the optimization component and the regret of the  
 257 bandit component. The former is simply the algorithm’s optimization cost minus the (overall) best

258 intervention’s optimization cost. The latter is defined similarly. Both of them need not be positive,  
259 but they sum up to the average regret of PROOF. This decomposition provides a rough illustration of  
260 how PROOF makes progress on both ends.

261 Under the experiment setting introduced above, Figure 2a shows that PROOF can quickly reduce the  
262 average regret in both optimization and bandit components. On the other hand, the performance of  
263 vanilla OFU is much more underwhelming. A typical bandit regret bound ignores many constant  
264 factors. We think this performance discrepancy is primarily due to the large variance in the implicit  
265 context  $x$  and  $c$ , which could be much better captured by training a predictive model. In fact, PROOF  
266 also has much smaller variance in its performance than vanilla OFU consistently.

267 We now tweak the problem parameters a bit and see how the performance changes. If we change the  
268 number of data points per iteration from  $n = 20$  to  $n = 40$ , we observe in Figure 2b that the regret of  
269 the optimization component becomes very small even at the beginning. This is because we have more  
270 data to learn from. If we change the upper bound of the norm of the linear mapping matrix  $F$  from  
271 10 to 100, we observe in Figure 2c that the optimization regret dominates the total regret. This is also  
272 reasonable because the optimization cost now has much larger magnitude than the bandit cost. Also,  
273 the vanilla OFU’s performance becomes quite disastrous, because now its “bandit” cost, which is the  
274 total cost for PROOF, has even larger magnitude and variance.

275 Admittedly, the setup above serves more as a proof of concept. In the second set of experiments, we  
276 scale up the experiments and show that PROOF still has better performance than the vanilla OFU  
277 baseline even when the problem parameters are not as friendly. Suppose we receive  $n = 500$  data  
278 points at every time step, and each data point has  $m = 50$  features. Keeping all other parameters  
279 the same, we observe in Fig. 2d that PROOF still significantly outperforms OFU. Note that in this  
280 case, we seem to have enough data points for the prediction task, and thus the bandit regret dominates  
281 the total regret. In Fig. 2e, we change the norm of the linear mapping matrix  $F$  from 10 to 1. This  
282 implies the optimization cost is now less important than the bandit cost. Indeed, this change reduces  
283 the variance for the OFU algorithm and thus it is doing much better than in previous experiments.  
284 However, our PROOF algorithm still outperforms OFU. Finally, in Fig. 2f, we increase the noise in  
285 the label distribution from  $\epsilon \sim \mathcal{N}(0, 0.1I_d)$  to  $\epsilon \sim \mathcal{N}(0, 0.5I_d)$ . This poses much more challenge to  
286 PROOF. As the data become more noisy, as expected, the optimization regret no longer stays close  
287 to zero as in the previous two experiments. Nevertheless, PROOF still manages to reduce the total  
288 regret at a faster rate than OFU.

## 289 **Declarations**

290 We believe that the positive social impact of our work has been thoroughly explained throughout this  
291 paper. We would like to stress here one potential negative impact that our particular work might bring.  
292 Our bandit data-driven optimization framework has an explicit structure to account for unintended  
293 social consequences of data-driven policy-making while the policy is being rolled out in the real  
294 world. This is our humble acknowledgment as AI researchers of the impossibility to consider all  
295 social factors in the initial formulation of an AI4SG problem. However, our work should not be  
296 misused as an excuse for not making the best effort to scope the unintended consequence ahead of  
297 time.

298 This work has not been published or first made available before January 1, 2017, and is original  
299 work.

## 300 **References**

- 301 [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear  
302 stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320,  
303 2011.
- 304 [2] Jacob Abernethy, Alex Chojnacki, Arya Farahi, Eric Schwartz, and Jared Webb. Activeremediation:  
305 The search for lead pipes in flint, michigan. In *Proceedings of the 24th ACM SIGKDD  
306 International Conference on Knowledge Discovery & Data Mining*, pages 5–14, 2018.
- 307 [3] Tom M Apostol. Introduction to analytic number theory. 1966.



- 308 [4] Susan Athey and Stefan Wager. Efficient policy learning. *arXiv preprint arXiv:1702.02896*,  
309 2017.
- 310 [5] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine*  
311 *Learning Research*, 3(Nov):397–422, 2002.
- 312 [6] Gah-Yi Ban and Cynthia Rudin. The big data newsvendor: Practical insights from machine  
313 learning. *Operations Research*, 67(1):90–108, 2019.
- 314 [7] Dimitris Bertsimas and Nathan Kallus. From predictive to prescriptive analytics. *Management*  
315 *Science*, 66(3):1025–1044, 2020.
- 316 [8] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff  
317 functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence*  
318 *and Statistics*, pages 208–214, 2011.
- 319 [9] Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under  
320 bandit feedback. In *Conference on Learning Theory*, 2008.
- 321 [10] Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning.  
322 *arXiv preprint arXiv:1103.4601*, 2011.
- 323 [11] Adam N Elmachtoub and Paul Grigas. "Smart" predict, then optimize". *arXiv preprint*  
324 *arXiv:1710.08005*, 2017.
- 325 [12] Fei Fang, Peter Stone, and Milind Tambe. When security games go green: designing defender  
326 strategies to prevent poaching and illegal fishing. In *Proceedings of the 24th International*  
327 *Conference on Artificial Intelligence*, pages 2589–2595, 2015.
- 328 [13] Nam Ho-Nguyen and Fatma Kılınc-Karzan. Risk guarantees for end-to-end prediction and  
329 optimization processes. 2019.
- 330 [14] Yi-hao Kao, Benjamin V Roy, and Xiang Yan. Directed regression. In *Advances in Neural*  
331 *Information Processing Systems*, pages 889–897, 2009.
- 332 [15] Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Ad-*  
333 *vances in applied mathematics*, 6(1):4–22, 1985.
- 334 [16] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to  
335 personalized news article recommendation. In *Proceedings of the 19th international conference*  
336 *on World wide web*, pages 661–670, 2010.
- 337 [17] Zheyuan Ryan Shi, Claire Wang, and Fei Fang. Artificial intelligence for social good: A survey.  
338 *arXiv preprint arXiv:2001.01818*, 2020.
- 339 [18] Zheyuan Ryan Shi, Yiwen Yuan, Kimberly Lo, Leah Lizarondo, and Fei Fang. Improving  
340 efficiency of volunteer-based food rescue operations. *Proceedings of the AAAI Conference on*  
341 *Artificial Intelligence*, 34(8):13369–13375, 2020.
- 342 [19] Adith Swaminathan and Thorsten Joachims. Batch learning from logged bandit feedback  
343 through counterfactual risk minimization. *The Journal of Machine Learning Research*,  
344 16(1):1731–1755, 2015.