

Learning to Interact with Learning Agents

Adish Singla
MPI-SWS
Saarbrücken, Germany
adishs@mpi-sws.org

Hamed Hassani
University of Pennsylvania
Philadelphia, USA
hassani@seas.upenn.edu

Andreas Krause
ETH Zurich
Zurich, Switzerland
krausea@ethz.ch

Abstract

AI and machine learning methods are increasingly interacting with and seeking information from people, robots, and other learning agents. Consequently, the learning dynamics of these agents creates fundamentally new challenges for existing methods. Motivated by the application of learning to offer personalized deals to users, we highlight these challenges by studying a variant of the framework of “online learning using expert advice with bandit feedback”. In our setting, we consider each expert as a learning agent, seeking to more accurately reflect real-world applications. The bandit feedback leads to additional challenges in this setting: at time t , only the expert i^t that has been selected by the central algorithm (forecaster) receives feedback from the environment and gets to learn at this time. A natural question to ask is whether it is possible to be competitive with the best expert j^* had it seen all the feedback, *i.e.*, competitive with the policy of always selecting expert j^* . We prove the following hardness result—without any coordination between the forecaster and the experts, it is impossible to design a forecaster achieving no-regret guarantees. We then consider a practical assumption allowing the forecaster to *guide* the learning process of the experts by blocking some of the feedback observed by them from the environment, *i.e.*, restricting the selected expert i^t to learn at time t for some time steps. With this additional coordination power, we design our forecaster LIL that achieves no-regret guarantees, and we provide regret bounds dependent on the learning dynamics of the best expert j^* .

Introduction

Many real-world applications involve repeatedly making decisions under uncertainty—for instance, choosing one of the several products to recommend to a user in an online recommendation service, or dynamically allocating resources among available stock options in a financial market. AI methods and machine learning techniques driving these applications are typically designed—or operate under the assumption—to be interacting with static components, *e.g.*, users’ preferences are fixed, or domain experts / trading tools providing stock recommendations are static. This assumption is often violated in modern applications as these methods are increasingly interacting with and seeking information from people, robots, and other learning agents. In this paper, we

highlight the fundamental challenges in designing algorithms that have to interact with learning agents, especially when algorithm’s decisions directly affect the learning dynamics of these agents.

Motivated by the application of learning to offer personalized deals to users, we focus on a well-studied framework of “online learning using expert advice with bandit feedback” (Littlestone and Warmuth 1994; Cesa-Bianchi et al. 1997; Freund and Schapire 1995; Auer et al. 2002; Cesa-Bianchi and Lugosi 2006; Bubeck and Cesa-Bianchi 2012). This is a generic framework for sequential decision making under uncertainty and addresses the fundamental question of how a learning algorithm should trade-off exploration (the cost of acquiring new information) versus exploitation (acting greedily based on current information to minimize instantaneous losses). In this paper, we investigate this framework with an important practical consideration: *How do we use the advice of experts when they themselves are learning agents?*

Motivating Applications

Modeling experts as learning agents realistically captures many practical scenarios of how one would define/encounter these experts in real-world applications, such as seeking advice from fellow players or friends, aggregating prediction recommendations from trading agents or different marketplaces, product testing with human participants who might adapt over time, and information acquisition from crowdsourcing participants who might learn over time. A specific instance of this problem setting is that of meta-learning whereby different learning algorithms (*e.g.*, with different hyperparameters or loss functions) are treated as experts (Baram, El-Yaniv, and Luz 2004; Hsu and Lin 2015; Maillard and Munos 2011; Agarwal et al. 2017).

As a concrete running example, we consider the problem of learning to offer personalized deals / discount coupons to users enabling new businesses to incentivize and attract more customers (Edelman, Jaffe, and Kominers 2011; Singla, Tschitschek, and Krause 2016; Hirnschall et al. 2018). An emerging trend is *deal-aggregator* sites like *Yipit*¹ providing personalized coupon recommendation services to their users by aggregating and selecting coupons from *daily-*

deal marketplaces like *Groupon* and *LivingSocial*¹. One of the primary goals of these recommendation systems like *Yipit* (corresponding to the central algorithm / forecaster in our setting) is to design better selection strategies for choosing coupons from different marketplaces (corresponding to the experts in our setting). However, these marketplaces themselves would be learning to optimize the coupons to offer, for instance, the discount price or the coupon type based on historic interactions with users (Edelman, Jaffe, and Kominers 2011).

Overview of Our Approach and Main Results

Our goal is to design a central online algorithm (henceforth, called as forecaster) to seek the advice of the available experts—more specifically, at time t , the forecaster selects an expert i^t , performs an action $a_{i^t}^t$ recommended by the expert i^t , and observes/incurs a loss $l^t(a_{i^t}^t)$ set by the adversary. Furthermore, given the bandit setting, only the selected expert i^t receives feedback from the environment and gets to learn at time t ; all other experts that have not been selected at time t experience no change in their learning state at this time. A natural benchmark in our problem setting is to be competitive with the best expert j^* had it seen all the feedback, *i.e.*, competitive with the cumulative loss one would incur by following the policy of always selecting expert j^* .

Generic setting and the hardness result. The fundamental challenge in our setting arises from the fact that the forecaster’s selection of experts affects which expert gets to learn at a particular time. In this paper, we establish the following hardness result—without any coordination between the forecaster and the experts, it is impossible to design a forecaster achieving no-regret guarantees when competing with the policy which always selects the expert j^* .

Additional coordination via blocking the feedback. In light of this hardness result, we next explore practically applicable approaches where it is possible to achieve no-regret for the forecaster. In order to make our results applicable to a wide range of real-world applications mentioned above, the focus of this paper is on a generic black-box approach in which the forecaster does not know the internal learning dynamics of the experts. The specific coordination protocol that we consider (alternatively, we can think of this as the additional power at the hands of the forecaster) is as follows: At a time t , the forecaster could decide to *block* the feedback from being observed by the selected expert i^t , thereby restricting the selected expert from learning at time t for some time steps. For instance, in the motivating application of offering personalized deals to users, the deal-aggregator site (forecaster) primarily interacts with users on behalf of the individual daily-deal marketplaces (experts) and hence could control the flow of feedback (*e.g.* users’ bids or clicks denoting their purchase decisions) to these marketplaces. With this additional coordination, we design our forecaster LIL that achieves no-regret guarantees with regret bounds dependent on the learning dynamics of the best expert j^* .

Connections to existing results. To conclude the overview of our results, we would like to point out a few relevant papers. First, Maillard and Munos (2011) introduced the EXP4/EXP3 algorithm, *i.e.*, EXP4 meta-

algorithm with experts executing EXP3 algorithms proving a regret bound of $\mathcal{O}(T^{\frac{2}{3}})$. Second, in a recent work contemporary to ours, Agarwal et al. (2017) provide improved regret bound of $\mathcal{O}(T^{\frac{1}{2}})$ (in comparison to the above-mentioned regret bound of $\mathcal{O}(T^{\frac{2}{3}})$) for the problem of designing meta-algorithm combining multiple bandit algorithms. Agarwal et al. (2017) also prove a hardness result similar in spirit to that of ours. However, all these existing meta-algorithms are based on the idea of feeding unbiased estimate of losses to the experts and are not directly applicable to our motivating applications where experts could be implementing learning algorithms with more complex feedback structure (*e.g.*, dynamic-pricing algorithm based on the partial monitoring framework (Cesa-Bianchi and Lugosi 2006; Bartók et al. 2014)), or experts being human agents who are learning over time. Below, we highlight two technical points of how our approach and main results differ from these existing results:

- Our coordination approach of blocking the feedback observed by experts (*i.e.*, making a binary decision, instead of modifying losses as done in existing approaches) is more suitable for real-world application scenarios, especially in situations where experts’ learning algorithms are not directly controlled by the forecaster and have complex feedback structure. Alternatively, when viewing this coordination in terms of communication between the forecaster and the selected expert i^t , our coordination can be achieved with a 1-bit of communication at time t , whereas the coordination in existing approaches requires communicating the probability of selecting expert i^t at time t .
- Our results apply to a rich class of no-regret online learning algorithms that experts might be implementing—the key ingredient of our results relies on proving a property, we termed as *smooth no-regret* learning dynamics. This property quantifies the robustness of an online learning algorithm w.r.t. the sparsity in the observed feedback and is of independent interest.

Generic Setting: The Model

In this section, we formally introduce the generic problem setting and discuss our objective. We have the following entities in our setting: (i) a central algorithm ALGO as the forecaster; (ii) an adversary ADV acting on behalf of the environment; and (iii) N experts $\text{EXP}_j \forall j \in \{1, \dots, N\}$ (henceforth denoted as $[N]$). Protocol 1 provides a high-level specification of the interaction between these entities. In subsequent sections, we will introduce an additional coordination allowing the forecaster to block the feedback observed by the selected expert for some time steps (*i.e.*, modifying the specification in line 7 of the Protocol 1).

Specification of the Interaction

The sequential decision making process proceeds in rounds $t = 1, 2, \dots, T$ (henceforth denoted as $[T]$); for simplicity we assume that T is known in advance to the forecaster and the results in this paper can be extended to an unknown horizon via the usual doubling trick (Cesa-Bianchi and Lugosi 2006).

Protocol 1: The interaction between adversary ADV, forecaster ALGO, and experts

```
foreach  $t = 1, 2, \dots, T$  do
  /* Adversary generates the following                                */
1  a private loss vector  $l^t$  for the forecaster, i.e.,  $l^t(a) \forall a \in \mathcal{A}$ 
2  a private feedback vector  $f^t$  for the experts, i.e.,  $f^t(a) \forall a \in \mathcal{A}$ 
  /* Selecting an expert and performing an action                        */
3  ALGO selects an expert  $i^t \in [N]$  denoted as  $\text{EXP}_{i^t}$ 
4  ALGO performs the action  $a_{i^t}^t$  recommended by  $\text{EXP}_{i^t}$ 
  /* Feedback and updates                                             */
5  ALGO incurs (and observes) loss  $l^t(a_{i^t}^t)$  and updates its selection strategy
6   $\forall j \in [N] : j \neq i^t$ ,  $\text{EXP}_j$  does not observe any feedback and makes no update
7   $\text{EXP}_{i^t}$  observes the feedback  $f^t(a_{i^t}^t)$  from the environment and updates its learning state
```

However, we do not assume that T is known to the experts. Each expert EXP_j where $j \in [N]$ is associated with a set of actions \mathcal{A}_j and the action set of the forecaster ALGO is given by $\mathcal{A} = \cup_{j \in [N]} \mathcal{A}_j$.² For the clarity of presentation in defining the loss and feedback vectors, we will consider that the action sets of experts are disjoint.³ The actions here could represent simple discrete actions (e.g., offering a discount coupon of a particular type) or could also represent functional policies defined over a time-dependent context (e.g., mapping user features to the value of a discount coupon).

At any time t , the adversary ADV generates a private loss vector l^t (i.e., $l^t(a) \forall a \in \mathcal{A}$) for the forecaster and a private feedback vector f^t (i.e., $f^t(a) \forall a \in \mathcal{A}$) for the experts—see examples below for the specific notion of feedback. Additionally, the adversary ADV generates a publicly available context that is accessible to all the experts at time t —this context essentially encodes any side information from the environment at time t (e.g., user’s features at time t).⁴

Simultaneously, the forecaster ALGO (possibly with some randomization) selects expert EXP_{i^t} to seek advice. The selected expert EXP_{i^t} recommends an action $a_{i^t} \in \mathcal{A}_{i^t} \subseteq \mathcal{A}$ (possibly with its internal randomization) which is then performed by the forecaster. The forecaster ALGO observes and incurs the loss $l^t(a_{i^t}^t)$, and updates its strategy on how to select experts in the future. All the experts apart from the one selected, $\text{EXP}_j \forall j \neq i^t$, observe no feedback and make no update—these experts do not experience *any* change in their learning state at this time. The selected expert EXP_{i^t} observes a feedback from the environment denoted as $f^t(a_{i^t}^t)$ and performs one learning step.

We assume that losses are bounded in the range $[0, l_{max}]$ for some known $l_{max} \in \mathbb{R}_+$; w.l.o.g. we will use $l_{max} = 1$ (Auer et al. 2002). We consider an oblivious (non-adaptive adversary) as is usual in the literature (Freund and Schapire 1995; Auer et al. 2002), i.e., the loss vector l^t and the feed-

back vector f^t at any time t do not depend on the actions taken by the forecaster, and hence can be considered to be fixed in advance. Apart from that, no other restrictions are put on the adversary, and it has complete knowledge about the forecaster and the learning dynamics of the experts.

The notion of the feedback and concrete examples. So far, we have considered a generic notion of the feedback received by the selected expert—this feedback essentially depends on the application setting and is supposed to be “compatible” with the learning algorithm used by an expert. For instance, consider an expert EXP_j implementing the EXP3 algorithm and recommending an action a_j^t at time t , then the feedback $f^t(a_j^t)$ received by this expert (if selected at time t) is the loss $l^t(a_j^t)$; for the case of expert EXP_j implementing the HEDGE algorithm, the feedback $f^t(a_j^t)$ received by this expert (if selected at time t) is the set of losses $\{l^t(a) \mid \forall a \in \mathcal{A}_j\}$. The feedback could be more general, for instance, receiving a binary signal of acceptance/rejection of the offered deal when an expert is implementing a dynamic-pricing algorithm based on the partial monitoring framework (Cesa-Bianchi and Lugosi 2006; Bartók et al. 2014).

Specification of the Experts

Next, we provide a formal specification of the experts. The focus of this paper is on a black-box approach in which the forecaster ALGO does not know the internal dynamics of the experts. At time t , let us denote an instance of feedback received by EXP_{i^t} by a tuple $h = (a_{i^t}^t, f^t(a_{i^t}^t))$. For any expert EXP_j where $j \in [N]$, let $\mathcal{H}_j^t = (h^1, h^2, \dots)$ denote the feedback history for EXP_j , i.e., an ordered sequence of feedback instances observed by EXP_j in the time period $[1, t)$. The length $|\mathcal{H}_j^t|$ denotes the number of learning steps for EXP_j up to time t . At time t , the action a_j^t recommended by EXP_j to the forecaster, if this expert is selected, is given by $a_j^t = \pi_j(\mathcal{H}_j^t)$ where π_j is a (possibly randomized) function of EXP_j , taking as input a history of feedback sequence, and outputs an action $a \in \mathcal{A}_j$. Importantly, this history \mathcal{H}_j^t is dependent on the execution of the forecaster ALGO—for clarity of presentation, we denote it as $\mathcal{H}_{j, \text{ALGO}}^t$.

No-regret learning dynamics. To be able to say anything meaningful in this setting, we introduce the constraint of *no-*

²The special case of standard multi-armed bandits is captured by the setting when \mathcal{A}_j is a singleton $\forall j \in [N]$.

³Note that assuming the disjoint action sets across experts is w.l.o.g., as we can still simulate the shared actions by enforcing that losses/feedbacks for the shared actions are same at any given time.

⁴We have omitted a formal definition of the context as it doesn’t directly play a role in the design of our forecaster.

regret learning dynamics on the experts. Let us consider any sequence of a loss vector l and a feedback vector f given by $\mathcal{D} = (l^\tau, f^\tau)_{\tau=\{1,2,\dots\}}$ generated arbitrarily by the adversary ADV and let $|\mathcal{D}|$ denotes its length. Consider a setting in which the forecaster executes a simple policy which always select a specific expert EXP_j for a fixed $j \in [N]$. Hence, this expert EXP_j gets to see all the feedback and has the *complete* feedback history at every time step—we denote this complete history at any time $\tau \in [|\mathcal{D}|]$ as $\mathcal{H}_{j,\text{FULL}}^\tau$. Then, the no-regret learning dynamics of EXP_j parameterized by $\beta_j \in [0, 1]$ guarantees that the cumulative loss of the forecaster executing this policy satisfies the following:

$$\mathbb{E} \left[\sum_{\tau=1}^{|\mathcal{D}|} l^\tau (\pi_j(\mathcal{H}_{j,\text{FULL}}^\tau)) \right] - \min_{a \in \mathcal{A}_j} \sum_{\tau=1}^{|\mathcal{D}|} l^\tau(a) \leq \mathcal{O}(|\mathcal{D}|^{\beta_j}) \quad (1)$$

where the expectation is w.r.t. the randomization of π_j .

Our Objective: No-Regret Guarantees

As a first attempt in designing the forecaster, one might consider using one of the standard algorithms from the EXP family (*e.g.*, the EXP3 algorithm (Auer et al. 2002) or the NEXP algorithm (McMahan and Streeter 2009)) as the forecaster. This would guarantee that the forecaster has no-regret guarantees using the classical notion of *external* regret (*cf.* Equation 2). However, we argue that external regret is not a desirable objective for our problem setting. We then formally state the guarantees we seek for the forecaster (*cf.* Equation 3).

External regret and its limitations. We begin by formally defining the classical notion of external regret used in the literature (Auer et al. 2002; Cesa-Bianchi and Lugosi 2006; Bubeck and Cesa-Bianchi 2012). Let us consider a complete execution of the forecaster ALGO in the retrospect: (i) let $\{a_j^t : j \in [N]\}_{t \in [T]}$ denote actions recommended by the experts during this execution and (ii) let $(l^t)_{t \in [T]}$ denote loss vectors generated by the adversary. In order to define the external regret of the forecaster in this execution, we need to *fix* these actions and loss vectors. Then, the external regret of the forecaster is given by

$$\text{REGEXT}(T, \text{ALGO}) := \sum_{t=1}^T l^t(a_{it}^t) - \min_{j \in [N]} \sum_{t=1}^T l^t(a_j^t) \quad (2)$$

If we would have used, let's say, the EXP3 algorithm (Auer et al. 2002) as the forecaster, we would obtain a bound of $\mathcal{O}(T^{\frac{1}{2}})$ on the external regret defined in Equation 2. However, this regret bound is only w.r.t. the post hoc sequence of actions performed and losses observed during the execution of the forecaster—it is not informative of the actual performance of the forecaster when comparing against a policy which always selects the best expert j^* (*cf.* Equation 3 for a formal definition). The reason why this classical notion of regret is not informative in our setting can also be attributed to the fact that losses at any time t are indirectly dependent on which experts were selected by the forecaster in the past as that defines the current

learning state of the experts (McMahan and Streeter 2009; Maillard and Munos 2011; Agarwal et al. 2017). This challenge of history-dependent losses also arises when playing against a non-oblivious/adaptive adversary and requires different notions of regret beyond the external regret (Arora, Dekel, and Tewari 2012).

Competing with the best expert. Intuitively, we want to be competitive with the best expert EXP_{j^*} (for any $j^* \in [N]$) had it seen all the feedback, *i.e.*, competitive with the policy of always selecting this expert EXP_{j^*} . In fact, executing such a policy ensures that the expert EXP_{j^*} gets full feedback to improve its learning state and perform well w.r.t. the single best action from the set \mathcal{A}_{j^*} . We can formally state this alternate notion of regret for the forecaster ALGO as follows:

$$\text{REG}(T, \text{ALGO}) := \sum_{t=1}^T \mathbb{E} \left[l^t (\pi_{it}(\mathcal{H}_{it, \text{ALGO}}^t)) \right] - \min_{j \in [N]} \min_{a \in \mathcal{A}_j} \sum_{t=1}^T l^t(a) \quad (3)$$

where the expectation is w.r.t. the randomization of the forecaster as well as any internal randomization of the experts.

If we already knew who the best expert EXP_{j^*} is at $t = 0$, we could always select this expert—the no-regret learning dynamics from Equation 1 dictates that the regret $\text{REG}(T, \text{ALGO})$ grows as $\mathcal{O}(T^{\beta_{j^*}})$. The main research question that we study in this paper is how to design an algorithm for the forecaster when we don't have this prior knowledge of j^* . It turns out that competing with this policy of always selecting EXP_{j^*} is a challenging problem in the bandit feedback setting (*cf.* next section for the hardness result). For instance, what might go wrong is that the best expert could have a slow rate of learning/convergence thus incurring high loss in the beginning, misleading the forecaster to essentially “downweigh” this expert. This in turn further exacerbates the problem for the best expert in the bandit feedback setting as this expert will be selected even less and thus have fewer learning steps to improve its state. This adds new challenges to the classic trade-off between exploration and exploitation, suggesting the need to explore at a higher rate.

Generic Setting: Hardness Result

In this section, we highlight the fundamental challenges in designing algorithms that have to interact with learning agents by establishing the following hardness result: in the absence of any coordination between the forecaster and the experts, it is impossible to design a forecaster that achieves no-regret guarantees (*cf.* Equation 3) in the worst-case. In fact, we prove this hardness result when playing against an oblivious (non-adaptive) adversary and when restricting the experts to be implementing well-behaved learning algorithms (*e.g.*, the HEDGE algorithm (Freund and Schapire 1995)). We formally state this hardness result in Theorem 1 below.

Theorem 1. *There is a setting in which each of the experts has no-regret learning dynamics with parameter $\beta \leq \frac{1}{2}$; however, any forecaster ALGO will suffer a linear regret, *i.e.*, $\text{REG}(T, \text{ALGO}) = \Omega(T)$.*

The proof is given in the extended version of this paper (Singla, Hassani, and Krause 2018); we briefly outline the

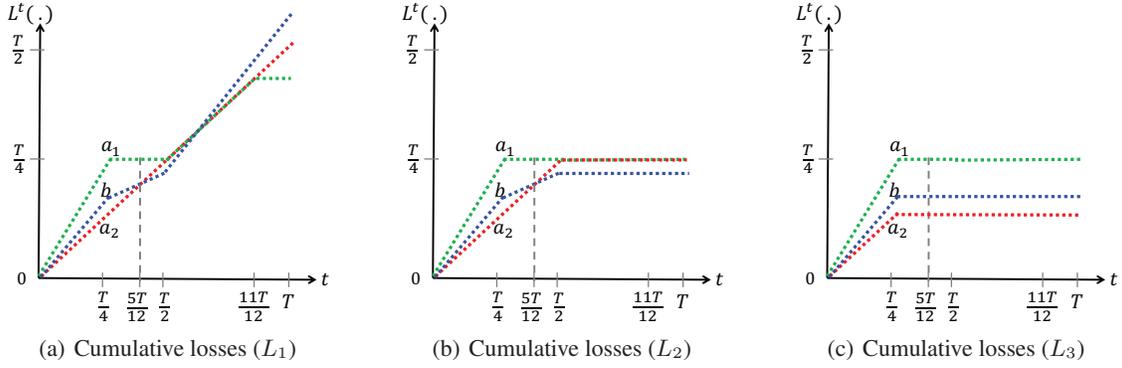


Figure 1: We have two experts: EXP_1 plays HEDGE and has two actions $\mathcal{A}_1 = \{a_1, a_2\}$, and EXP_2 has only one action $\mathcal{A}_2 = \{b\}$. Figures 1(a), 1(b), and 1(c) shows the *cumulative* loss sequences L_1, L_2 , and L_3 for three different scenarios—the adversary at $t = 0$ uniformly at random picks one of these scenarios and uses that loss sequence. These plots show the cumulative losses of the three actions $\mathcal{A} = \{a_1, a_2, b\}$ for three different sequences. The losses are illustrated with the following color scheme— a_1 :green, a_2 :red, and b :blue.

main ideas below. Our setting for proving this theorem consists of two experts EXP_1 and EXP_2 . The first expert EXP_1 has two actions given by $\mathcal{A}_1 = \{a_1, a_2\}$, and the second expert EXP_2 has only one action given by $\mathcal{A}_2 = \{b\}$. The action set of the forecaster ALGO is given by $\mathcal{A} = \{a_1, a_2, b\}$. The expert EXP_1 plays the HEDGE algorithm (Freund and Schapire 1995), *i.e.*, the regret rate parameter is $\beta_1 = 0.5$; the expert EXP_2 has only one action to play, hence $\beta_2 = 0$. Figures 1(a), 1(b), and 1(c) show the *cumulative* loss sequences L_1, L_2 , and L_3 for three different scenarios—the adversary at $t = 0$ uniformly at random picks one of these scenarios and uses that loss sequence.

Our main argument is that for any forecaster, one of the scenarios leads to linear regret.⁵ In the proof, we consider the case where a forecaster is facing the sequence L_1 . We then divide the time horizon T into different slots, and discuss the execution behavior of the forecaster and experts over these time slots. Specifically, our claim is that in the time slot $t \in (\frac{T}{4}, \frac{T}{2}]$, the expert EXP_1 would not be selected for $\frac{T}{12} - o(T)$ time steps. As a result, in the time slot $t \in (\frac{11T}{12}, T]$, the expert EXP_1 would end up recommending action a_2 , and a_1 would only be recommended $o(T)$ number of times, leading to $\Omega(T)$ regret for the forecaster. Informally speaking, our negative example shows that the forecaster’s selection strategy could add “blind spots” in the feedback history of the experts and that they might not be able to “recover” from this. The fundamental challenge leading to this hardness result is that the forecaster’s selection strategy affects the feedback observed by the experts, which in turn alters the learning processes of these experts.

⁵In fact, this hardness result holds even when considering a powerful forecaster which knows exactly the learning algorithms used by the experts, and is able to see the losses $\{l^t(a_1), l^t(a_2), l^t(b)\}$ at every time $t \in [T]$.

Our Approach

In this section, we explore practically applicable approaches where it is possible to achieve no-regret guarantees for the forecaster. Our setting is similar to the problem of designing meta-algorithm for combining different bandit algorithms (Maillard and Munos 2011; Bubeck and Cesa-Bianchi 2012; Agarwal et al. 2017). However, existing meta-algorithms—for instance, the EXP4/EXP3 algorithm (Maillard and Munos 2011) or the CORRAL algorithm (Agarwal et al. 2017)—are based on the idea that the forecaster has the power to modify losses as seen by the experts, thereby, feeding unbiased estimate of losses to the experts. However, these existing meta-algorithms are not directly applicable to our motivating applications where experts could be implementing learning algorithms with more complex feedback structure (*e.g.*, dynamic-pricing algorithm based on the partial monitoring framework (Cesa-Bianchi and Lugosi 2006; Bartók et al. 2014)), or experts being human learning agents.

Additional Coordination and Forecaster LIL

Motivated by the application setting of deal-aggregator sites, we consider a new coordination approach in which the forecaster has the power to guide the experts’ learning process by carefully blocking the feedback observed by them. For instance, in the motivating application of offering personalized deals to users, the deal-aggregator site (forecaster) primarily interacts with users on behalf of the individual daily-deal marketplaces (experts) and hence could control the flow of feedback (*e.g.* users’ bids or clicks denoting their purchase decisions) to these marketplaces. Formally, at a time t , the forecaster could decide to block the feedback from being observed by the selected expert i^t , thereby restricting the selected expert from learning at time t for some time steps (*i.e.*, modifying the specification in line 7 of the Protocol 1). Compared to the existing meta-algorithms which require modifying losses, our coordination approach is more applicable to generic black-box settings where experts’ learning

Algorithm 2: Forecaster LIL

```
1 Parameters:  $\eta \in (0, 1]$ 
2 Initialize: time  $t = 1$ , weights  $w_j^t = 1 \forall j \in [N]$ 
   foreach  $t = 1, 2, \dots, T$  do
   | /* Selecting an expert */
   |  $\forall j \in [N]$ , define probability
   |  $p_j^t = (1 - \eta) \cdot \frac{w_j^t}{(\sum_{k \in [N]} w_k^t)} + \frac{\eta}{N}$ 
   | Draw  $i^t$  from the multinomial distribution
   |  $(p_j^t)_{j \in [N]}$ 
   | Perform action  $a_{i^t}^t$  recommended by  $\text{EXP}_{i^t}$ 
   | /* Making updates */
   | Observe loss  $l^t(a_{i^t}^t)$ 
   |  $\forall j \in [N]$ , do the following:
   |   Set  $\tilde{l}_j^t = \frac{l^t(a_{i^t}^t)}{p_{i^t}^t}$  for  $j = i^t$ , else  $\tilde{l}_j^t = 0$ 
   |   Update  $w_j^{t+1} \leftarrow w_j^t \cdot \exp(-\frac{\eta \cdot \tilde{l}_j^t}{N})$ 
   | /* Blocking the feedback */
   |  $\xi^t \sim \text{Bernoulli}(\frac{\eta}{N \cdot p_{i^t}^t})$ 
   | if ( $\xi^t = 0$ ) then
   |   EXP $_{i^t}$  does not observe the feedback  $f^t(a_{i^t}^t)$ 
   |   and has no change in the learning state
```

algorithms are not directly controlled by the forecaster and could have different feedback structure.

With this additional coordination power, we design our forecaster LIL (**L**earning to **I**nteract with **L**earners), presented in Algorithm 2. The selection strategy of the forecaster LIL is similar to the EXP3 algorithm (Auer et al. 2002). The core idea of guiding the experts’ learning process is presented in lines 10, 11, and 12 of the Algorithm 2. At time t , the forecaster LIL blocks the feedback observed by the selected expert when $\xi^t = 0$ where ξ^t is a $\text{Bernoulli}(\frac{\eta}{N \cdot p_{i^t}^t})$ random variable. By choosing this particular random variable, the forecaster LIL ensures that the probability that *any* expert $\text{EXP}_j \forall j \in [N]$ observes feedback is constant over time $t \in [T]$ and is given by η/N . Considering the negative example used in the proof of Theorem 1, this means that by carefully restricting the feedback, the forecaster LIL avoids any “blind spots” in the feedback history of the experts. However, in order to achieve this, LIL is required to explore at a higher rate, as is evident by the value of η in Theorem 2.

Performance Analysis

In this section, we analyze the performance of the forecaster LIL. We first introduce a novel property, we termed as *smooth no-regret* learning dynamics, quantifying the robustness of an online learning algorithm w.r.t. the sparsity in the observed feedback. Then, the key ingredient of our results (cf. Theorem 2) relies on showing that a rich class of no-regret online learning algorithms that experts might be implementing also satisfy this property of smoothness (cf. Proposition 1). This property is of independent interest towards designing learning

algorithms that are robust against deletion of feedback.

Theoretical Guarantees

Smooth no-regret learning dynamics. In the bandit feedback setting, not all the experts get to observe feedback at a given time, and hence the feedback history would not be *complete* for the experts (cf. Equation 1). Consider the same setting as used in defining the no-regret learning dynamics in Equation 1 for a specific expert EXP_j . Again, the forecaster executes a simple policy which always select this expert EXP_j , however, let us consider a situation where the expert EXP_j only gets to observe the feedback with a probability $\alpha \in (0, 1]$ at a given time. We denote this *sparse* feedback history at any time $\tau \in [D]$ as $\mathcal{H}_{j, \alpha\text{-FULL}}^\tau$. Then, the *smooth no-regret* learning dynamics of EXP_j guarantees that the cumulative loss of the forecaster satisfies the following:

$$\mathbb{E} \left[\sum_{\tau=1}^{|\mathcal{D}|} l^\tau(\pi_j(\mathcal{H}_{j, \alpha\text{-FULL}}^\tau)) \right] - \min_{a \in \mathcal{A}_j} \sum_{\tau=1}^{|\mathcal{D}|} l^\tau(a) \leq \mathcal{O} \left(\frac{(\alpha \cdot |\mathcal{D}|)^{\beta_j}}{\alpha} \right) \quad (4)$$

where the expectation is w.r.t. the randomization of function π_j as well as w.r.t. the randomization in generating this sparse feedback history. The following proposition states that a rich class of online learning algorithms indeed have smooth no-regret learning dynamics that can be used by the experts—the proof is given in the extended version of this paper (Singla, Hassani, and Krause 2018).

Proposition 1. *A rich class of no-regret online learning algorithms based on gradient-descent style updates have smooth learning dynamics including the Online Mirror Descent family of algorithms with exact or estimated gradients (Shalev-Shwartz 2011) and Online Convex Programming via greedy projections (Zinkevich 2003).*

No-regret guarantees of LIL. Next, we prove the no-regret guarantees of our forecaster LIL when competing with the best expert EXP_{j^*} (cf. Equation 3). The following theorem states the bounds, keeping only the leading terms of T . The proof is given in the extended version of this paper (Singla, Hassani, and Krause 2018).

Theorem 2. *Let T be the fixed time horizon. Consider that the best expert $j^* \in [N]$ has no-regret smooth learning dynamics parameterized by $\beta_{j^*} \in [0, 1]$ and LIL is invoked with input $\beta \in [0, 1]$ such that $\beta \geq \beta_{j^*}$. Set parameters $\eta = \Theta(T^{-\frac{1-\beta}{2-\beta}} \cdot N^{\frac{1-\beta}{2-\beta}} \cdot (\log N)^{\frac{1}{2} \cdot \mathbf{1}_{\{\beta=0\}}})$. Then, for sufficiently large T , the worst-case expected cumulative regret of the forecaster LIL is:*

$$\text{REG}(T, \text{LIL}) \leq \mathcal{O}(T^{\frac{1}{2-\beta}} \cdot N^{\frac{1}{2-\beta}} \cdot (\log N)^{\frac{1}{2} \cdot \mathbf{1}_{\{\beta=0\}}})$$

For the special case of multi-armed bandits (where $\beta = 0$), this regret bound matches the bound of $\Theta(T^{\frac{1}{2}})$ —in fact, for this special case, our algorithm LIL is exactly equivalent to EXP3. For an important case when experts are implementing algorithms like HEDGE or EXP3 (where $\beta = \frac{1}{2}$), our algorithm LIL achieves the bound of $\mathcal{O}(T^{\frac{2}{3}})$.

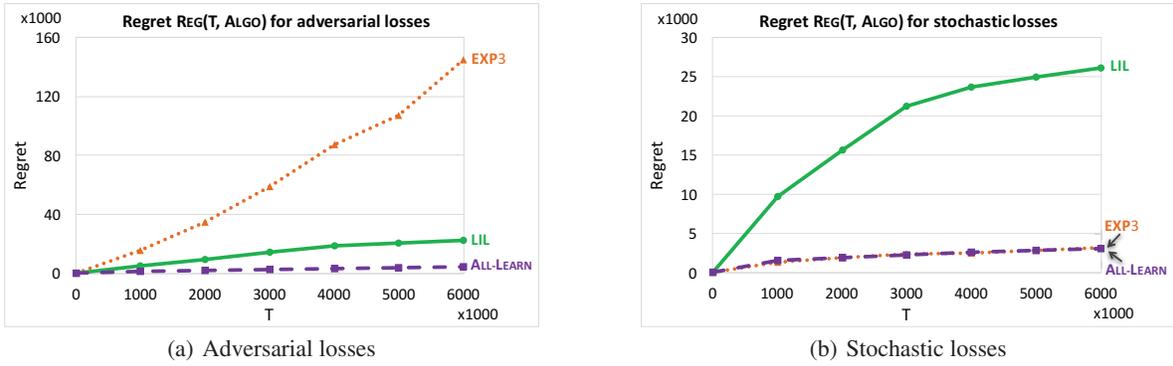


Figure 2: Simulation results showing the performance of LIL against different benchmarks for a setting with two experts (similar to what we used in the proof of Theorem 1). (a) considers adversarial losses from Figure 1(a); (b) considers stochastic losses. Results illustrate that LIL has provable no-regret guarantees in general, and also a fast convergence rate in the stochastic setting.

Simulation Results

Next, we evaluate the performance of the forecaster LIL via simulations, and compare against the following benchmarks:

- **EXP3**: using EXP3 algorithm (Auer et al. 2002) as the forecaster for the specification in Protocol 1.
- **ALL-LEARN**: using EXP3 algorithm (Auer et al. 2002) as the forecaster for a relaxed/easier setting in which all experts $j \in [N]$ observe the feedback at any time t .

Adversarial losses. As our first simulation setting, we consider the same set up used in the proof of Theorem 1 and we use the loss sequence shown in Figure 1(a). For this loss sequence, the loss of actions $A = \{a_1, a_2, b\}$ averaged over $t \in [T]$ is given by $(0.4583, 0.5, 0.7487)$ —hence the best expert is EXP₁ and the best action is a_1 (cf. Equation 3). Figure 2(a) shows the regret $\text{REG}(T, \text{ALGO})$ for LIL, EXP3, and ALL-LEARN, and illustrates the following points. First, EXP3 suffer a linear regret, as dictated by the hardness result in Theorem 1. Second, LIL has a sub-linear regret as proved in Theorem 2. Note that if we plot $\text{REG}(T, \text{ALGO})$ shown in Figure 2(a) on a $\log\text{-}\log$ plot, the slope s of a linear fit on the resulting plot defines the rate T^s of the growth of regret. The slope is $s = 0.62$ for LIL—our results from Theorem 2 dictate an upper bound of 0.66 on LIL’s slope (for $\beta = 0.5$).

Stochastic losses. To complement the above results, we next consider a stochastic version of the above setup where losses of actions $A = \{a_1, a_2, b\}$ are sampled i.i.d. from Bernoulli distributions with means given by $(0.45, 0.5, 0.475)$ —as before, the best expert is EXP₁ and the best action is a_1 for this stochastic setting. Figure 2(b) shows the regret $\text{REG}(T, \text{ALGO})$ for LIL, EXP3, and ALL-LEARN, and illustrates the following points. First, EXP3 performs better than LIL: this is expected because in the stochastic setting, the strategy to block the feedback only slows down convergence; furthermore, LIL has a higher rate of exploration η compared to EXP3. Second, in the $\log\text{-}\log$ plot, the slope is $s = 0.56$ for LIL, $s = 0.47$ for EXP3, and $s = 0.40$ for ALL-LEARN—this signifies the fast convergence rate of LIL. Third, the regret of LIL in Figure 2(a) (adversarial) and Figure 2(b) (stochastic) is about the same because LIL’s

coordination approach essentially adds stochasticity in the feedback observed by experts.

Further Related Work

Markovian, rested, and restless bandits. Our setting is similar in spirit to that of the *rested* Markovian bandits where each action/arm is associated with its own stochastic MDP and an arm changes its state only when it is pulled. In the seminal work, Gittins (1979) introduced the *Gittins index* to find an optimal sequential policy for these Markovian bandits problem. This work has been extended to settings where all arms change their reward distributions at every time step according to their associated stochastic MDPs, termed as *restless* bandits (Whittle 1988; Slivkins and Upfal 2008; Besbes, Gur, and Zeevi 2014). However, none of these frameworks could model learning dynamics of the experts in the adversarial setting we consider.

Online boosting and adaptive control. Another line of recent work similar in spirit to ours is online boosting—combining a set of “weak” online learning algorithms to form a “strong” online learning algorithm (Beygelzimer, Kale, and Luo 2015; Beygelzimer et al. 2015). However, there are substantial differences when compared to our problem setting: online boosting techniques have been studied in the context of classification/regression problems, (most of) the techniques are for the full-information setting, and one of the key challenges revolves around defining the weights for an online training example to be passed on to the “weak” learners. Our work is also similar to that of *switching adaptive control* (Fu 2015) which employs an array of simple candidate controllers; the goal is to design a meta-controller that can switch across controllers to search for the best candidate controller in real-time. Our result in Theorem 1 highlights the challenges in designing a meta-controller in an adversarial setting.

Learning in games. An orthogonal line of research studies the interaction of agents in multiplayer games where each agent uses a no-regret learning algorithm (Blum and Mansour 2007; Syrgkanis et al. 2015). The questions tackled in this line of research are very different as it focuses on interactions of the agents, their individual as well as social

utilities, and the convergence of the game to an equilibrium. This orthogonal line of research reassures that the no-regret learning dynamics that we consider in this paper are indeed important and natural dynamics that are also prevalent in other application domains.

Conclusions and Future Work

In this paper, we investigated the framework of online learning using expert advice with bandit feedback when experts themselves are learning agents. Our hardness result highlights the fundamental challenges faced by traditional AI and machine learning methods when interacting with learning agents. In order to circumvent the hardness result, we introduced a new coordination approach allowing the forecaster to guide the experts' learning process by restricting the feedback received by them. In comparison to existing meta-algorithms that modify losses seen by the experts, our approach is more suitable for real-world applications where learning agents might have more complex feedback structure—for instance, a deal-aggregator site interacting with daily-deal marketplaces, or an AI system interacting with humans.

An important direction would be to study other practical ways of coordination and to understand the minimal communication required between the central algorithm and learning agents to achieve desired guarantees. For our problem setting, an interesting question to tackle is whether it is possible to design a forecaster using our coordination approach (which requires only 1-bit of communication at every time step) with a cumulative regret of $\Theta(T^{\frac{1}{2}})$ when the individual experts have no-regret learning dynamics with parameter $\beta = \frac{1}{2}$. Finally, we would like to point out that while we focused on the framework of online learning using expert advice, our results call for further studies of other frameworks and methods, e.g., active learning methods when dealing with dynamic oracles.

Acknowledgments. This work was supported in part by the Swiss National Science Foundation, Nano-Tera.ch program as part of the Opensense II project, ERC StG 307036, and a Microsoft Research Faculty Fellowship. Adish Singla acknowledges support by a Facebook Graduate Fellowship.

References

- Agarwal, A.; Luo, H.; Neyshabur, B.; and Schapire, R. E. 2017. Corraling a band of bandit algorithms. In *COLT*, 12–38.
- Arora, R.; Dekel, O.; and Tewari, A. 2012. Online bandit learning against an adaptive adversary: from regret to policy regret. In *ICML*.
- Auer, P.; Cesa-Bianchi, N.; Freund, Y.; and Schapire, R. E. 2002. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing* 32(1):48–77.
- Baram, Y.; El-Yaniv, R.; and Luz, K. 2004. Online choice of active learning algorithms. *Journal of Machine Learning Research*.
- Bartók, G.; Foster, D. P.; Pál, D.; Rakhlin, A.; and Szepesvári, C. 2014. Partial monitoring – Classification, regret bounds, and algorithms. *Mathematics of Operations Research*.
- Besbes, O.; Gur, Y.; and Zeevi, A. J. 2014. Optimal exploration-exploitation in a multi-armed-bandit problem with non-stationary rewards. In *NIPS*.
- Beygelzimer, A.; Hazan, E.; Kale, S.; and Luo, H. 2015. Online gradient boosting. In *NIPS*.
- Beygelzimer, A.; Kale, S.; and Luo, H. 2015. Optimal and adaptive algorithms for online boosting. In *ICML*, 2323–2331.
- Blum, A., and Mansour, Y. 2007. Learning, regret minimization, and equilibria.
- Bubeck, S., and Cesa-Bianchi, N. 2012. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Machine Learning* 5(1):1–122.
- Cesa-Bianchi, N., and Lugosi, G. 2006. *Prediction, learning, and games*. Cambridge university press.
- Cesa-Bianchi, N.; Freund, Y.; Haussler, D.; Helmbold, D. P.; Schapire, R. E.; and Warmuth, M. K. 1997. How to use expert advice. *Journal of the ACM* 44(2):427–485.
- Edelman, B.; Jaffe, S.; and Kominers, S. D. 2011. Togroupon or not togroupon: The profitability of deep discounts. *Marketing Letters* 1–15.
- Freund, Y., and Schapire, R. E. 1995. A decision-theoretic generalization of on-line learning and an application to boosting. In *COLT*, 23–37.
- Fu, M. 2015. Switching adaptive control. In *Encyclopedia of Systems and Control*. Springer.
- Gittins, J. C. 1979. Bandit processes and dynamic allocation indices. *Journal of the Royal Statistical Society. Series B (Methodological)* 148–177.
- Hirschall, C.; Singla, A.; Tschitschek, S.; and Krause, A. 2018. Learning user preferences to incentivize exploration in the sharing economy. In *AAAI*.
- Hsu, W.-N., and Lin, H.-T. 2015. Active learning by learning. In *AAAI*, 2659–2665.
- Littlestone, N., and Warmuth, M. K. 1994. The weighted majority algorithm. *Info and Computation* 70(2):212–261.
- Maillard, O.-A., and Munos, R. 2011. Adaptive bandits: Towards the best history-dependent strategy. In *AISTATS*.
- McMahan, H. B., and Streeter, M. J. 2009. Tighter bounds for multi-armed bandits with expert advice. In *COLT*.
- Shalev-Shwartz, S. 2011. Online learning and online convex optimization. *Foundations and Trends in Machine Learning* 4(2):107–194.
- Singla, A.; Hassani, H.; and Krause, A. 2018. Learning to interact with learning agents (extended version).
- Singla, A.; Tschitschek, S.; and Krause, A. 2016. Actively learning hemimetrics with applications to eliciting user preferences. In *ICML*.
- Slivkins, A., and Upfal, E. 2008. Adapting to a changing environment: the brownian restless bandits. In *COLT*.
- Syrkkanis, V.; Agarwal, A.; Luo, H.; and Schapire, R. E. 2015. Fast convergence of regularized learning in games. In *NIPS*, 2971–2979.
- Whittle, P. 1988. Restless bandits: Activity allocation in a changing world. *Journal of applied probability*.
- Zinkevich, M. 2003. Online convex programming and generalized infinitesimal gradient ascent. In *ICML*.