# Nonstochastic Multi-Armed Bandits
# with Graph-Structured Feedback *

Noga Alon
Tel Aviv University, Tel Aviv, Israel,
Institute for Advanced Study, Princeton, United States,
and Microsoft Research, Herzliya, Israel
nogaa@tau.ac.il

Nicolò Cesa-Bianchi
Università degli Studi di Milano, Italy
nicolo.cesa-bianchi@unimi.it

Claudio Gentile
University of Insubria, Italy
claudio.gentile@uninsubria.it

Shie Mannor
Technion, Israel
shie@ee.technion.ac.il

Yishay Mansour
Microsoft Research and
Tel-Aviv University, Israel
mansour@tau.ac.il

Ohad Shamir
Weizmann Institute of Science, Israel
ohad.shamir@weizmann.ac.il

**Abstract**

We present and study a partial-information model of online learning, where a decision maker repeatedly chooses from a finite set of actions, and observes some subset of the associated losses. This naturally models several situations where the losses of different actions are related, and knowing the loss of one action provides information on the loss of other actions. Moreover, it generalizes and interpolates between the well studied full-information setting (where all losses are revealed) and the bandit setting (where only the loss of the action chosen by the player is revealed). We provide several algorithms addressing different variants of our setting, and provide tight regret bounds depending on combinatorial properties of the information feedback structure.

## 1  Introduction

Prediction with expert advice —see, e.g., [8, 9, 15, 19, 23]— is a general abstract framework for studying sequential decision problems. For example, consider a weather forecasting problem, where each day we receive predictions from various experts, and we need to devise our forecast. At the end of the day, we observe how well each expert did, and we can use this information to improve our forecasting in the future. Our goal is that over time, our performance converges to that of

---

*Preliminary versions of this manuscript appeared in [1, 20].

the best expert in hindsight. More formally, such problems are often modeled as a repeated game between a player and an adversary, where each round, the adversary privately assigns a loss value to each action in a fixed set (in the example above, the discrepancy in the forecast if we follows a given expert's advice). Then the player chooses an action (possibly using randomization), and incurs the corresponding loss. The goal of the player is to control regret, which is defined as the cumulative excess loss incurred by the player as compared to the best fixed action over a sequence of rounds.

In some situations, however, the player only gets partial feedback on the loss associated with each action. For example, consider a web advertising problem, where every day one can choose an ad to display to a user, out of a fixed set of ads. As in the forecasting problem, we sequentially choose actions from a given set, and may wish to control our regret with respect to the best fixed ad in hindsight. However, while we can observe whether a displayed ad was clicked on, we do not know what would have happened if we chose a different ad to display. In our abstract framework, this corresponds to the player observing the loss of the action picked, but not the losses of other actions. This well-known setting is referred to as the (non-stochastic) multi-armed bandit problem, which in this paper we denote as the *bandit* setting. In contrast, we refer to the previous setting, where the player observes the losses of all actions, as the *expert* setting. In this work, our main goal is to bridge between these two feedback settings, and create a spectrum of models in between.

Before continuing, let us first quantify the performance attainable in the expert and the bandit setting. Letting $K$ be the number of available actions, and $T$ be the number of played rounds, the best possible regret for the expert setting is of order $\sqrt{\ln(K)\,T}$. This optimal rate is achieved by the Hedge algorithm [15] or the Follow the Perturbed Leader algorithm [17]. In the bandit setting, the optimal regret is of order $\sqrt{KT}$, achieved by the INF algorithm [3]. A bandit variant of Hedge, called Exp3 [4], achieves a regret with a slightly worse bound of order $\sqrt{K\ln(K)\,T}$. Thus, switching from the full-information expert setting to the partial-information bandit setting increases the attainable regret by a multiplicative factor of $\sqrt{K}$, up to extra logarithmic factors. This exponential difference in terms of the dependence on $K$ can be crucial in problems with large action sets. The intuition for this difference in performance has long been that in the bandit setting, we only get $1/K$ of the information obtained in the expert setting (as we observe just a single loss, rather than all $K$ at each round), hence the additional $K$-factor under the square root in the bound.

While the bandit setting received much interest, it can be criticized for not capturing additional side-information we often have on the losses of the different actions. As a motivating example, consider the problem of web advertising mentioned earlier. In the standard multi-armed bandits setting, we assume that we have no information whatsoever on whether undisplayed ads would have been clicked on. However, in many relevant cases, the semantic relationship among actions (ads) implies that we do indeed have some side-information. For instance, if two ads $i$ and $j$ are for similar vacation packages in Hawaii, and ad $i$ was displayed and clicked on by some user, it is likely that the other ad $j$ would have been clicked on as well. In contrast, if ad $i$ is for high-end running shoes, and ad $j$ is for wheelchair accessories, then a user who clicked on one ad is unlikely to click on the other. This sort of side-information is not captured by the standard bandit setting. A similar type of side-information arises in product recommendation systems hosted on online social networks, in which users can befriend each other. In this case, it has been observed that social relationships reveal similarities in tastes and interests [21]. Hence, a product liked by some user may also be liked by the user's friends. A further example, not in the marketing domain, is route selection: We are given a graph of possible routes connecting cities. When we select a route connecting two

cities, we observe the cost (say, driving time or fuel consumption) of the "edges" along that route and, in addition, we have complete information on sub-routes including any subset of the edges.[1]

In this paper, we present and study a setting which captures these types of side-information, and in fact interpolates between the bandit setting and the expert setting. This is done by defining a *feedback system*, under which choosing a given action also reveals the losses of some subset of the other actions. This feedback system can be viewed as a directed and time-changing graph $G_t$ over actions: an arc (directed edge) from action $i$ to action $j$ implies that when playing action $i$ at round $t$ we get information also about the loss of action $j$ at round $t$. Thus, the expert setting is obtained by choosing a complete graph over actions (playing any action reveals all losses), and the bandit setting is obtained by choosing an empty edge set (playing an action only reveals the loss of that action). The attainable regret turns out to depend on non-trivial combinatorial properties of this graph. To describe our results, we need to make some distinctions in the setting that we consider.

**Directed vs. symmetric setting.**   In some situations, the side-information between two actions is symmetric —for example, if we know that both actions will have a similar loss. In that case, we can model our feedback system $G_t$ as an undirected graph. In contrast, there are situations where the side-information is not symmetric. For example, consider the side-information gained from asymmetric social links, such as followers of celebrities. In such cases, followers might be more likely to shape their preferences after the person they follow, than the other way around. Hence, a product liked by a celebrity is probably also liked by his/her followers, whereas a preference expressed by a follower is more often specific to that person. Another example in the context of ad placement is when a person buying a video game console might also buy a high-def cable to connect it to the TV set. Vice versa, interest in high-def cables need not indicate an interest in game consoles. In such situations, modeling the feedback system via a directed graph $G_t$ is more suitable. Note that the symmetric setting is a special case of the directed setting, and therefore handling the symmetric case is easier than the directed case.

**Informed vs. uninformed setting.**   In some cases, the feedback system is known to the player before each round, and can be utilized for choosing actions. For example, we may know beforehand which pairs of ads are related, or we may know the users who are friends of another user. We denote this setting as the informed setting. In contrast, there might be cases where the player does not have full knowledge of the feedback system before choosing an action, and we denote this harder setting as the uninformed setting. For example, consider a firm recommending products to users of an online social network. If the network is owned by a third party, and therefore not fully visible, the system may still be able to run its recommendation policy by only accessing small portions of the social graph around each chosen action (i.e., around each user to whom a recommendation is sent).

Generally speaking, our contribution lies in both characterizing the regret bounds that can be achieved in the above settings as a function of combinatorial properties of the feedback systems, as well as providing efficient sequential decision algorithms working in those settings. More specifically,

---

[1] Though this example may also be viewed as an instance of combinatorial bandits [10], the model we propose is more general. For example, it does not assume linear losses, which could arise in the routing example from the partial ordering of sub-routes.

our contributions can be summarized as follows (see Section 2 for a brief review of the relevant combinatorial properties of graphs).

**Uninformed setting.** We present an algorithm (Exp3-SET) that achieves $\widetilde{\mathcal{O}}\left(\sqrt{\ln(K)\sum_{t=1}^{T}\mathtt{mas}(G_t)}\right)$ regret in expectation, where $\mathtt{mas}(G_t)$ is the size of the maximal acyclic graph in $G_t$. In the symmetric setting, $\mathtt{mas}(G_t) = \alpha(G_t)$ ($\alpha(G_t)$ is the independence number of $G_t$), and we prove that the resulting regret bound is optimal up to logarithmic factors, when $G_t = G$ is fixed for all rounds. Moreover, we show that Exp3-SET attains $\mathcal{O}\left(\sqrt{\ln(K)\,T}\right)$ regret when the feedback graphs $G_t$ are random graphs generated from a standard Erdős-Renyi model.

**Informed setting.** We present an algorithm (Exp3-DOM) that achieves expected regret of $\mathcal{O}\left(\ln(K)\sqrt{\ln(KT)\sum_{t=1}^{T}\alpha(G_t)}\right)$, for both the symmetric and directed cases. Since our lower bound also applies to the informed setting, this characterizes the attainable regret in the informed setting, up to logarithmic factors. Moreover, we present another algorithm (ELP.P), that achieves $\mathcal{O}\left(\sqrt{\ln(K/\delta)\sum_{t=1}^{T}\mathtt{mas}(G_t)}\right)$ regret with probability at least $1 - \delta$ over the algorithm's internal randomness. Such a high-probability guarantee is stronger than the guarantee for Exp3-DOM, which holds just in expectation, and turns out to be of the same order in the symmetric case. However, in the directed case, the regret bound may be weaker since $\mathtt{mas}(G_t)$ may be larger than $\alpha(G_t)$. Moreover, ELP.P requires us to solve a linear program at each round, whereas Exp3-DOM only requires finding an approximately minimal dominating set, which can be done by a standard greedy set cover algorithm.

Our results interpolate between the bandit and expert settings: When $G_t$ is a full graph for all $t$ (which means that the player always gets to see all losses, as in the expert setting), then $\mathtt{mas}(G_t) = \alpha(G_t) = 1$, and we recover the standard guarantees for the expert setting: $\sqrt{T}$ up to logarithmic factors. In contrast, when $G_t$ is the empty graph for all $t$ (which means that the player only observes the loss of the action played, as in the bandit setting), then $\mathtt{mas}(G_t) = \alpha(G_t) = K$, and we recover the standard $\sqrt{KT}$ guarantees for the bandit setting, up to logarithmic factors. In between are regret bounds scaling like $\sqrt{BT}$, where $B$ lies between 1 and $K$, depending on the graph structure (again, up to log-factors).

Our results are based on the algorithmic framework for handling the standard bandit setting introduced in [4]. In this framework, the full-information Hedge algorithm is combined with unbiased estimates of the full loss vectors in each round. The key challenge is designing an appropriate randomized scheme for choosing actions, which correctly balances exploration and exploitation or, more specifically, ensures small regret while simultaneously controlling the variance of the loss estimates. In our setting, this variance is subtly intertwined with the structure of the feedback system. For example, a key quantity emerging in the analysis of Exp3-DOM can be upper bounded in terms of the independence number of the graphs. This bound (Lemma 16 in the appendix) is based on a combinatorial construction which may be of independent interest.

For the uninformed setting, our work was recently improved by [18], whose main contribution is an algorithm attaining $\mathcal{O}\left(\sqrt{\ln(K)\ln(KT)\sum_{t=1}^{T}\alpha(G_t)}\right)$ expected regret in the uninformed and directed setting using a novel implicit exploration idea. Up to log factors, this matches the performance of our Exp3-DOM and ELP.P algorithms, without requiring prior knowledge of the

feedback system. On the other hand, their bound holds only in expectation rather than with high probability.

**Paper Organization:** In the next section, we formally define our learning protocols, introduce our main notation, and recall the combinatorial properties of graphs that we require. In Section 3, we tackle the uninformed setting, by introducing Exp3-SET, with upper and lower bounds on regret based on both the size of the maximal acyclic subgraph (general directed case) and the independence number (symmetric case). In Section 4, we handle the informed setting through the two algorithms Exp3-DOM (Section 4.1) on which we prove regret bounds in expectation, and ELP.P (Section 4.2) whose bounds hold in the more demanding high probability regime. We conclude the main text with Section 5, where we discuss open questions, and possible directions for future research. All technical proofs are provided in the appendices. We organized such proofs based on which section of the main text the corresponding theoretical claims occur.

## 2 Learning protocol, notation, and preliminaries

As stated in the introduction, we consider adversarial decision problems with a finite action set $V = \{1, \ldots, K\}$. At each time $t = 1, 2, \ldots$, a player (the "learning algorithm") picks some action $I_t \in V$ and incurs a bounded loss $\ell_{I_t,t} \in [0, 1]$. Unlike the adversarial bandit problem [4, 9], where only the played action $I_t$ reveals its loss $\ell_{I_t,t}$, here we assume all the losses in a subset $S_{I_t,t} \subseteq V$ of actions are revealed after $I_t$ is played. More formally, the player observes the pairs $(i, \ell_{i,t})$ for each $i \in S_{I_t,t}$. We also assume $i \in S_{i,t}$ for any $i$ and $t$, that is, any action reveals its own loss when played. Note that the bandit setting ($S_{i,t} = \{i\}$) and the expert setting ($S_{i,t} = V$) are both special cases of this framework. We call $S_{i,t}$ the *feedback set* of action $i$ at time $t$, and write $i \xrightarrow{t} j$ when at time $t$ playing action $i$ also reveals the loss of action $j$. (We sometimes write $i \rightarrow j$ when time $t$ plays no role in the surrounding context.) With this notation, $S_{i,t} = \{j \in V : i \xrightarrow{t} j\}$. The family of feedback sets $\{S_{i,t}\}_{i \in V}$ we collectively call the *feedback system* at time $t$.

The adversaries we consider are nonoblivious. Namely, each loss $\ell_{i,t}$ and feedback set $S_{i,t}$ at time $t$ can be arbitrary functions of the past player's actions $I_1, \ldots, I_{t-1}$ (note, though, that the regret is measured with respect to a fixed action assuming the adversary would have chosen the same losses, so our results do not extend to truly adaptive adversaries in the sense of [13]). The performance of a player $A$ is measured through the expected regret

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big]$$

where $L_{A,T} = \ell_{I_1,1} + \cdots + \ell_{I_T,T}$ and $L_{k,T} = \ell_{k,1} + \cdots + \ell_{k,T}$ are the cumulative losses of the player and of action $k$, respectively.[2] The expectation is taken with respect to the player's internal randomization (since losses are allowed to depend on the player's past random actions, $L_{k,T}$ may also be random). In Section 3 we also consider a variant in which the feedback system is randomly generated according to a specific stochastic model. For simplicity, we focus on a finite horizon setting, where the number of rounds $T$ is known in advance. This can be easily relaxed using a standard doubling trick.

---

[2] Although we defined the problem in terms of losses, our analysis can be applied to the case when actions return rewards $g_{i,t} \in [0, 1]$ via the transformation $\ell_{i,t} = 1 - g_{i,t}$.

We also consider the harder setting where the goal is to bound the actual regret

$$L_{A,T} - \max_{k \in V} L_{k,T}$$

with high probability $1 - \delta$ with respect to the player's internal randomization, and where the regret bound depends logarithmically on $1/\delta$. Clearly, a high probability bound on the actual regret implies a similar bound on the expected regret.

Whereas some of our algorithms need to know the feedback system at the beginning of each step $t$, others need it only at the end of each step. We thus consider two online learning settings: the *informed* setting, where the full feedback system $\{S_{i,t}\}_{i \in V}$ selected by the adversary is made available to the learner *before* making the choice $I_t$; and the *uninformed setting*, where no information whatsoever regarding the time-$t$ feedback system is given to the learner prior to prediction, but only following the prediction and with the associated loss information.

We find it convenient at this point to adopt a graph-theoretic interpretation of feedback systems. At each step $t = 1, 2, \ldots, T$, the feedback system $\{S_{i,t}\}_{i \in V}$ defines a directed graph $G_t = (V, D_t)$, the feedback graph, where $V$ is the set of actions and $D_t$ is the set of arcs (i.e., ordered pairs of nodes). For $j \neq i$, the arc $(i, j)$ belongs to $D_t$ if and only if $i \xrightarrow{t} j$ (the self-loops created by $i \xrightarrow{t} i$ are intentionally ignored). Hence, we can equivalently define $\{S_{i,t}\}_{i \in V}$ in terms of $G_t$. Observe that the outdegree $d_{i,t}^+$ of any $i \in V$ equals $|S_{i,t}| - 1$. Similarly, the indegree $d_{i,t}^-$ of $i$ is the number of actions $j \neq i$ such that $i \in S_{j,t}$ (i.e., such that $j \xrightarrow{t} i$). A notable special case of the above is when the feedback system is symmetric: $j \in S_{i,t}$ if and only if $i \in S_{j,t}$ for all $i, j$ and $t$. In words, playing $i$ at time $t$ reveals the loss of $j$ if and only if playing $j$ at time $t$ reveals the loss of $i$. A symmetric feedback system defines an undirected graph $G_t$ or, more precisely, a directed graph having, for every pair of nodes $i, j \in V$, either no arcs or length-two directed cycles. Thus, from the point of view of the symmetry of the feedback system, we also distinguish between the *directed* case ($G_t$ is a general directed graph) and the *symmetric* case ($G_t$ is an undirected graph for all $t$).

The analysis of our algorithms depends on certain properties of the sequence of graphs $G_t$. Two graph-theoretic notions playing an important role here are those of *independent sets* and *dominating sets*. Given an undirected graph $G = (V, E)$, an independent set of $G$ is any subset $T \subseteq V$ such that no two $i, j \in T$ are connected by an edge in $E$, i.e., $(i, j) \notin E$. An independent set is *maximal* if no proper superset thereof is itself an independent set. The size of any largest (and thus maximal) independent set is the *independence number* of $G$, denoted by $\alpha(G)$. If $G$ is directed, we can still associate with it an independence number: we simply view $G$ as undirected by ignoring arc orientation. If $G = (V, D)$ is a directed graph, then a subset $R \subseteq V$ is a dominating set for $G$ if for all $j \notin R$ there exists some $i \in R$ such that $(i, j) \in D$. In our bandit setting, a time-$t$ dominating set $R_t$ is a subset of actions with the property that the loss of any remaining action in round $t$ can be observed by playing some action in $R_t$. A dominating set is *minimal* if no proper subset thereof is itself a dominating set. The domination number of directed graph $G$, denoted by $\gamma(G)$, is the size of any smallest (and therefore minimal) dominating set for $G$; see Figure 1 for examples.

Computing a minimum dominating set for an arbitrary directed graph $G_t$ is equivalent to solving a minimum set cover problem on the associated feedback system $\{S_{i,t}\}_{i \in V}$. Although minimum set cover is NP-hard, the well-known Greedy Set Cover algorithm [12], which repeatedly selects from $\{S_{i,t}\}_{i \in V}$ the set containing the largest number of uncovered elements so far, computes a dominating set $R_t$ such that $|R_t| \leq \gamma(G_t)(1 + \ln K)$.

We can also lift the notion of independence number of an undirected graph to directed graphs through the notion of *maximum acyclic subgraphs*. Given a directed graph $G = (V, D)$, an acyclic
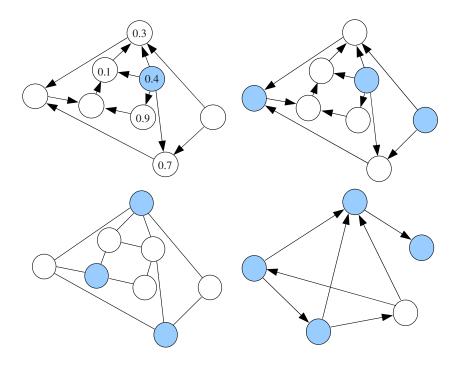
6

Figure 1: An example for some graph-theoretic concepts. **Top Left:** A feedback system with $K = 8$ actions (self-loops omitted). The light blue action reveals its loss 0.4, as well as the losses of the other four actions it points to. **Top Right:** The light blue nodes are a minimal dominating set for the same graph. The rightmost action is included in any dominating set, since no other action is dominating it. **Bottom Left:** A symmetric feedback system where the light blue nodes are a maximal independent set. This is the same graph as before, but edge orientation has been removed. **Bottom Right:** The light blue nodes are a maximum acyclic subgraph of the depicted 5-action graph.

subgraph of $G$ is any graph $G' = (V', D')$ such that $V' \subseteq V$, and $D' = D \cap (V' \times V')$, with no (directed) cycles. We denote by $\mathtt{mas}(G) = |V'|$ the maximum size of such $V'$. Note that when $G$ is undirected (more precisely, as above, when $G$ is a directed graph having for every pair of nodes $i, j \in V$ either no arcs or length-two cycles), then $\mathtt{mas}(G) = \alpha(G)$, otherwise $\mathtt{mas}(G) \geq \alpha(G)$. In particular, when $G$ is itself a directed acyclic graph, then $\mathtt{mas}(G) = |V|$. See Figure 1 (bottom right) for a simple example. Finally, we let $\mathbb{I}\{A\}$ denote the indicator function of event $A$.

## 3 The uninformed setting

In this section we investigate the setting in which the learner must select an action without any knowledge of the current feedback system. We introduce a simple general algorithm, Exp3-SET (Algorithm 1), that works in both the directed and symmetric cases. In the symmetric case, we show that the regret bound achieved by the algorithm is optimal to within logarithmic factors.

When the feedback graph $G_t$ is a fixed clique or a fixed edgeless graph, Exp3-SET reduces to the Hedge algorithm or, respectively, to the Exp3 algorithm. Correspondingly, the regret bound

---

**Algorithm 1:** The Exp3-SET algorithm (for the uninformed setting)

**Parameter:** $\eta \in [0, 1]$

**Initialize:** $w_{i,1} = 1$ for all $i \in V = \{1, \ldots, K\}$

**For** $t = 1, 2, \ldots$ :

1. Feedback system $\{S_{i,t}\}_{i \in V}$ and losses $\ell_t$ are generated but not disclosed ;

2. Set $p_{i,t} = \dfrac{w_{i,t}}{W_t}$ for each $i \in V$, where $W_t = \displaystyle\sum_{j \in V} w_{j,t}$ ;

3. Play action $I_t$ drawn according to distribution $p_t = (p_{1,t}, \ldots, p_{K,t})$ ;

4. Observe:

   (a) pairs $(i, \ell_{i,t})$ for all $i \in S_{I_t, t}$;

   (b) Feedback system $\{S_{i,t}\}_{i \in V}$ is disclosed;

5. For any $i \in V$ set $w_{i,t+1} = w_{i,t} \exp(-\eta \, \widehat{\ell}_{i,t})$, where

$$\widehat{\ell}_{i,t} = \frac{\ell_{i,t}}{q_{i,t}} \mathbb{I}\{i \in S_{I_t, t}\} \qquad \text{and} \qquad q_{i,t} = \sum_{j : j \xrightarrow{t} i} p_{j,t} \ .$$

---

for Exp3-SET yields the regret bound of Hedge and that of Exp3 as special cases.

Similar to Exp3, Exp3-SET uses importance sampling loss estimates $\widehat{\ell}_{i,t}$ that divide each observed loss $\ell_{i,t}$ by the probability $q_{i,t}$ of observing it. This probability $q_{i,t}$ is the probability of observing the loss of action $i$ at time $t$, i.e., it is simply the sum of all $p_{j,t}$ (the probability of selecting action $j$ at time $t$) such that $j \xrightarrow{t} i$ (recall that this sum always includes $p_{i,t}$).

In the expert setting, we have $q_{i,t} = 1$ for all $i$ and $t$, and we recover the Hedge algorithm. In the bandit setting, $q_{i,t} = p_{i,t}$ for all $i$ and $t$, and we recover the Exp3 algorithm (more precisely, we recover the variant Exp3Light of Exp3 that does not have an explicit exploration term, see [11] and also [22, Theorem 2.7]).

In what follows, we show that the regret of Exp3-SET can be bounded in terms of the key quantity

$$Q_t = \sum_{i \in V} \frac{p_{i,t}}{q_{i,t}} = \sum_{i \in V} \frac{p_{i,t}}{\sum_{j : j \xrightarrow{t} i} p_{j,t}} \ . \tag{1}$$

Each term $p_{i,t}/q_{i,t}$ can be viewed as the probability of drawing $i$ from $p_t$ conditioned on the event that $\ell_{i,t}$ was observed. A key aspect to our analysis is the ability to deterministically and non-vacuously[3] upper bound $Q_t$ in terms of certain quantities defined on $\{S_{i,t}\}_{i \in V}$. We do so in two ways, either irrespective of how small each $p_{i,t}$ may be (this section) or depending on suitable lower bounds on the probabilities $p_{i,t}$ (Section 4). In fact, forcing lower bounds on $p_{i,t}$ is equivalent to

---

[3] An obvious upper bound on $Q_t$ is $K$, since $p_{i,t}/q_{i,t} \leq 1$.

adding exploration terms to the algorithm, which can be done only when $\{S_{i,t}\}_{i \in V}$ is known before each prediction (i.e., in the informed setting).

The following result, whose proof is in Appendix A.2, is the building block for all subsequent results in the uninformed setting.

**Lemma 1** *The regret of Exp3-SET satisfies*

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \mathbb{E}[Q_t] \; . \tag{2}$$

In the expert setting, $q_{i,t} = 1$ for all $i$ and $t$ implies $Q_t = 1$ deterministically for all $t$. Hence, the right-hand side of (2) becomes $(\ln K)/\eta + (\eta/2)\,T$ , corresponding to the Hedge bound with a slightly larger constant in the second term; see, e.g., [9, Page 72]. In the bandit setting, $q_{i,t} = p_{i,t}$ for all $i$ and $t$ implies $Q_t = K$ deterministically for all $t$. Hence, the right-hand side of (2) takes the form $(\ln K)/\eta + (\eta/2)\,KT$ , equivalent to the Exp3 bound; see, e.g., [5, Equation 3.4].

We now move on to the case of general feedback systems, for which we can prove the following result (proof is in Appendix A.3).

**Theorem 2** *The regret of Exp3-SET satisfies*

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \mathbb{E}[\mathtt{mas}(G_t)] \; .$$

*If $\mathtt{mas}(G_t) \leq m_t$ for $t = 1, \ldots, T$, then setting $\eta = \sqrt{(2 \ln K)\big/ \sum_{t=1}^{T} m_t}$ gives*

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq \sqrt{2(\ln K) \sum_{t=1}^{T} m_t} \; .$$

As we pointed out in Section 2, $\mathtt{mas}(G_t) \geq \alpha(G_t)$, with equality holding when $G_t$ is an undirected graph. Hence, in the special case when $G_t$ is symmetric, we obtain the following result.

**Corollary 3** *In the symmetric case, the regret of Exp3-SET satisfies*

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \mathbb{E}[\alpha(G_t)] \; .$$

*If $\alpha(G_t) \leq \alpha_t$ for $t = 1, \ldots, T$, then setting $\eta = \sqrt{(2 \ln K)\big/ \sum_{t=1}^{T} \alpha_t}$ gives*

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq \sqrt{2(\ln K) \sum_{t=1}^{T} \alpha_t} \; .$$

Note that both Theorem 2 and Corollary 3 require the algorithm to know upper bounds on $\mathtt{mas}(G_t)$ and $\alpha(G_t)$, which may be computationally non-trivial – we return and expand on this issue in section 4.2.

In light of Corollary 3, one may wonder whether Lemma 1 is powerful enough to allow a control of regret in terms of the independence number even in the directed case. Unfortunately, the next result shows that —in the directed case— $Q_t$ cannot be controlled unless specific properties of $p_t$ are assumed. More precisely, we show that even for simple directed graphs, there exist distributions $p_t$ on the vertices such that $Q_t$ is linear in the number of nodes while the independence number[4] is 1.

**Fact 4** *Let $G = (V, D)$ be a total order on $V = \{1, \ldots, K\}$, i.e., such that for all $i \in V$, arc $(j, i) \in D$ for all $j = i + 1, \ldots, K$. Let $p = (p_1, \ldots, p_K)$ be a distribution on $V$ such that $p_i = 2^{-i}$, for $i < K$ and $p_k = 2^{-K+1}$. Then*

$$Q = \sum_{i=1}^{K} \frac{p_i}{p_i + \sum_{j : j \to i} p_j} = \sum_{i=1}^{K} \frac{p_i}{\sum_{j=i}^{K} p_j} = \frac{K + 1}{2} \ .$$

Next, we discuss lower bounds on the achievable regret for arbitrary algorithms. The following theorem provides a lower bound on the regret in terms of the independence number $\alpha(G)$, for a constant graph $G_t = G$ (which may be directed or undirected).

**Theorem 5** *Suppose $G_t = G$ for all $t$ with $\alpha(G) > 1$. There exist two constants $C_1, C_2 > 0$ such that whenever $T \geq C_1 \alpha(G)^3$, then for any algorithm there exists an adversarial strategy for which the expected regret of the algorithm is at least $C_2 \sqrt{\alpha(G)T}$.*

The intuition of the proof (provided in Appendix A.4) is the following: if the graph $G$ has $\alpha(G)$ non-adjacent vertices, then an adversary can make this problem as hard as a standard bandit problem, played on $\alpha(G)$ actions. Since for bandits on $K$ actions there is a $\Omega(\sqrt{KT})$ lower bound on the expected regret, a variant of the proof technique leads to a $\Omega(\sqrt{\alpha(G)T})$ lower bound in our case.

One may wonder whether a sharper lower bound exists which applies to the general directed adversarial setting and involves the larger quantity $\mathtt{mas}(G)$. Unfortunately, the above measure does not seem to be related to the optimal regret: using Lemma 11 in Appendix A.5 (see proof of Theorem 6 below) one can exhibit a sequence of graphs each having a large acyclic subgraph, on which the regret of Exp3-SET is still small.

**Random feedback systems.** We close this section with a study of Lemma 1 in a setting where the feedback system is stochastically generated via the Erdős-Renyi model. This is a standard model for random directed graphs $G = (V, D)$, where we are given a density parameter $r \in [0, 1]$ and, for any pair $i, j \in V$, arc $(i, j) \in D$ with independent probability $r$ (self loops, i.e., arcs $(i, i)$ are included by default here). We have the following result.

**Theorem 6** *For $t = 1, 2, \ldots$, let $G_t$ be an independent draw from the Erdős-Renyi model with fixed parameter $r \in [0, 1]$. Then the regret of Exp3-SET satisfies*

$$\max_{k \in V} \mathbb{E}[L_{A,T} - L_{k,T}] \leq \frac{\ln K}{\eta} + \frac{\eta T}{2r}\left(1 - (1 - r)^K\right) \ .$$

---

[4] In this specific example, the maximum acyclic subgraph has size $K$, which confirms the looseness of Theorem 2.

*In the above, expectations are computed with respect to both the algorithm's randomization and the random generation of $G_t$ occurring at each round. In particular, setting $\eta = \sqrt{\frac{2r \ln K}{T\left(1-(1-r)^K\right)}}$, gives*

$$\max_{k \in V} \mathbb{E}\left[L_{A,T} - L_{k,T}\right] \le \sqrt{\frac{2(\ln K)T\left(1-(1-r)^K\right)}{r}} \;.$$

Note that as $r$ ranges in $[0, 1]$ we interpolate between the multi-arm bandit[5] ($r = 0$) and the expert ($r = 1$) regret bounds.

Finally, note that standard results from the theory of Erdős-Renyi graphs —at least in the symmetric case (see, e.g., [16])— show that when the density parameter $r$ is constant, the independence number of the resulting graph has an inverse dependence on $r$. This fact, combined with the lower bound above, gives a lower bound of the form $\sqrt{T/r}$, matching (up to logarithmic factors) the upper bound of Theorem 6.

# 4    The informed setting

The lack of a lower bound matching the upper bound provided by Theorem 2 is a good indication that something more sophisticated has to be done in order to upper bound the key quantity $Q_t$ defined in (1). This leads us to consider more refined ways of allocating probabilities $p_{i,t}$ to nodes. We do so by taking advantage of the informed setting, in which the learner can access $G_t$ before selecting the action $I_t$. The algorithm Exp3-DOM, introduced in this section, exploits the knowledge of $G_t$ in order to achieve an optimal (up to logarithmic factors) regret bound.

Recall the problem uncovered by Fact 4: when the graph induced by the feedback system is directed, $Q_t$ cannot be upper bounded, in a non-vacuous way, independent of the choice of probabilities $p_{i,t}$. The new algorithm Exp3-DOM controls these probabilities by adding an exploration term to the distribution $p_t$. This exploration term is supported on a dominating set of the current graph $G_t$, and computing such a dominating set before selection of the action at time $t$ can only be done in the informed setting. Intuitively, exploration on a dominating set allows to control $Q_t$ by increasing the probability $q_{i,t}$ that each action $i$ is observed. If the dominating set is also minimal, then the variance caused by exploration can be bounded in terms of the independence number (and additional logarithmic factors) just like the undirected case.

Yet another reason why we may need to know the feedback system beforehand is when proving high probability results on the regret. In this case, operating with a feedback term for the probabilities $p_{i,t}$ seems unavoidable. In Section 4.2 we present another algorithm, called ELP.P, which can deliver regret bounds that hold with high probability over its internal randomization.

## 4.1    Bounds in expectation: the Exp3-DOM algorithm

The Exp3-DOM algorithm (see Algorithm 2) for the informed setting runs $\mathcal{O}(\log K)$ variants of Exp3 (with explicit exploration) indexed by $b = 0, 1, \ldots, \lfloor \log_2 K \rfloor$. At time $t$ the algorithm is given the current feedback system $\{S_{i,t}\}_{i \in V}$, and computes a dominating set $R_t$ of the directed graph $G_t$ induced by $\{S_{i,t}\}_{i \in V}$. Based on the size $|R_t|$ of $R_t$, the algorithm uses instance $b_t = \lfloor \log_2 |R_t| \rfloor$ to

---

[5] Observe that $\lim_{r \to 0^+} \frac{1-(1-r)^K}{r} = K$.

---

**Algorithm 2:** The Exp3-DOM algorithm (for the informed setting)

---

**Input:** Exploration parameters $\gamma^{(b)} \in (0, 1]$ for $b \in \{0, 1, \dots, \lfloor \log_2 K \rfloor\}$

**Initialization:** $w_{i,1}^{(b)} = 1$ for all $i \in V = \{1, \dots, K\}$ and $b \in \{0, 1, \dots, \lfloor \log_2 K \rfloor\}$

**For** $t = 1, 2, \dots$ :

1. Feedback system $\{S_{i,t}\}_{i \in V}$ is generated *and disclosed*, (losses $\ell_t$ are generated and not disclosed);

2. Compute a dominating set $R_t \subseteq V$ for $G_t$ associated with $\{S_{i,t}\}_{i \in V}$ ;

3. Let $b_t$ be such that $|R_t| \in [2^{b_t}, 2^{b_t+1} - 1]$;

4. Set $W_t^{(b_t)} = \sum_{i \in V} w_{i,t}^{(b_t)}$;

5. Set $p_{i,t}^{(b_t)} = \left(1 - \gamma^{(b_t)}\right) \dfrac{w_{i,t}^{(b_t)}}{W_t^{(b_t)}} + \dfrac{\gamma^{(b_t)}}{|R_t|} \mathbb{I}\{i \in R_t\}$;

6. Play action $I_t$ drawn according to distribution $p_t^{(b_t)} = \left(p_{1,t}^{(b_t)}, \dots, p_{K,t}^{(b_t)}\right)$ ;

7. Observe pairs $(i, \ell_{i,t})$ for all $i \in S_{I_t,t}$;

8. For any $i \in V$ set $w_{i,t+1}^{(b_t)} = w_{i,t}^{(b_t)} \exp\left(-\gamma^{(b_t)} \widehat{\ell}_{i,t}^{(b_t)} / 2^{b_t}\right)$, where

$$\widehat{\ell}_{i,t}^{(b_t)} = \frac{\ell_{i,t}}{q_{i,t}^{(b_t)}} \mathbb{I}\{i \in S_{I_t,t}\} \qquad \text{and} \qquad q_{i,t}^{(b_t)} = \sum_{j \,:\, j \xrightarrow{t} i} p_{j,t}^{(b_t)} \ .$$

---

draw action $I_t$. We use a superscript $b$ to denote the quantities relevant to the variant of Exp3 indexed by $b$. Similarly to the analysis of Exp3-SET, the key quantities are

$$q_{i,t}^{(b)} = \sum_{j \,:\, i \in S_{j,t}} p_{j,t}^{(b)} = \sum_{j \,:\, j \xrightarrow{t} i} p_{j,t}^{(b)} \qquad \text{and} \qquad Q_t^{(b)} = \sum_{i \in V} \frac{p_{i,t}^{(b)}}{q_{i,t}^{(b)}} \ , \qquad b = 0, 1, \dots, \lfloor \log_2 K \rfloor \ .$$

Let $T^{(b)} = \{t = 1, \dots, T \,:\, |R_t| \in [2^b, 2^{b+1} - 1]\}$. Clearly, the sets $T^{(b)}$ are a partition of the time steps $\{1, \dots, T\}$, so that $\sum_b |T^{(b)}| = T$. Since the adversary adaptively chooses the dominating sets $R_t$ (through the adaptive choice of the feedback system at time $t$), the sets $T^{(b)}$ are random variables. This causes a problem in tuning the parameters $\gamma^{(b)}$. For this reason, we do not prove a regret bound directly for Exp3-DOM, where each instance uses a fixed $\gamma^{(b)}$, but for a slight variant of it (described in the proof of Lemma 7 — see Appendix B.1), where each $\gamma^{(b)}$ is set through a doubling trick.

**Lemma 7** *In the directed case, the regret of Exp3-DOM satisfies*

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big] \le \sum_{b=0}^{\lfloor \log_2 K \rfloor} \left( \frac{2^b \ln K}{\gamma^{(b)}} + \gamma^{(b)} \mathbb{E}\left[ \sum_{t \in T^{(b)}} \left( 1 + \frac{Q_t^{(b)}}{2^{b+1}} \right) \right] \right) . \tag{3}$$

*Moreover, if we use a doubling trick to choose $\gamma^{(b)}$ for each $b = 0, \ldots, \lfloor \log_2 K \rfloor$, then*

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big] = \mathcal{O}\left( (\ln K) \, \mathbb{E}\left[ \sqrt{\sum_{t=1}^{T} \left( 4|R_t| + Q_t^{(b_t)} \right)} \right] + (\ln K) \ln(KT) \right) . \tag{4}$$

Importantly, the next result (proof in Appendix B.2) shows how bound (4) of Lemma 7 can be expressed in terms of the sequence $\alpha(G_t)$ of independence numbers of graphs $G_t$ whenever the Greedy Set Cover algorithm [12] (see Section 2) is used to compute the dominating set $R_t$ of the feedback system at time $t$.

**Theorem 8** *If Step 2 of Exp3-DOM uses the Greedy Set Cover algorithm to compute the dominating sets $R_t$, then the regret of Exp-DOM using the doubling trick satisfies*

$$\max_{k \in V} \mathbb{E}\big[L_{A,T} - L_{k,T}\big] = \mathcal{O}\left( \ln(K) \sqrt{\ln(KT) \sum_{t=1}^{T} \alpha(G_t)} + \ln(K) \ln(KT) \right) .$$

Combining the upper bound of Theorem 8 with the lower bound of Theorem 5, we see that the attainable expected regret in the informed setting is characterized by the independence numbers of the graphs. Moreover, a quick comparison between Corollary 3 and Theorem 8 reveals that a symmetric feedback system overcomes the advantage of working in an informed setting: The bound we obtained for the uninformed symmetric setting (Corollary 3) is sharper by logarithmic factors than the one we derived for the informed — but more general, i.e., directed — setting (Theorem 8).

## 4.2 High probability bounds: the ELP.P algorithm

We now turn to present an algorithm working in the informed setting for which we can also prove high-probability regret bounds.[6] We call this algorithm ELP.P (which stands for "Exponentially-weighted algorithm with Linear Programming", with high Probability). Like Exp3-DOM, the exploration component is not uniform over the actions, but is chosen carefully to reflect the graph structure at each round. In fact, the optimal choice of the exploration for ELP.P requires us to solve a simple linear program, hence the name of the algorithm.[7] The pseudo-code appears as Algorithm 3. Note that unlike the previous algorithms, this algorithm utilizes the "rewards" formulation of the problem, i.e., instead of using the losses $\ell_{i,t}$ directly, it uses the rewards $g_{i,t} = 1 - \ell_{i,t}$, and boosts the weight of actions for which $g_{i,t}$ is estimated to be large, as opposed to decreasing the weight of actions for which $\ell_{i,t}$ is estimated to be large. This is done merely for technical convenience, and does not affect the complexity of the algorithm nor the regret guarantee.

---

[6] We have been unable to prove high-probability bounds for Exp3-DOM or variants of it.

[7] We note that this algorithm improves over the basic ELP algorithm initially presented in [20], in that its regret is bounded in high probability and not just in expectation, and applies in the directed case as well as the symmetric case.

---
**Algorithm 3:** The ELP.P algorithm (for the informed setting)
---

**Input:** Confidence parameter $\delta \in (0,1)$, learning rate $\eta > 0$;

**Initialization:** $w_{i,1} = 1$ for all $i \in V = \{1, \ldots, K\}$;

**For** $t = 1, 2, \ldots$ :

1. Feedback system $\{S_{i,t}\}_{i \in V}$ is generated *and disclosed*, (losses $\ell_t$ are generated and not disclosed);

2. Let $\Delta_K$ be the $K$-dimensional probability simplex, and $s_t = (s_{1,t}, \ldots s_{K,t})$ be a solution to the linear program
$$\max_{(s_1, \ldots, s_K) \in \Delta_K} \quad \min_{i \in V} \sum_{j : j \xrightarrow{t} i} s_j$$

3. Set $p_{i,t} := (1 - \gamma_t) \frac{w_{i,t}}{W_t} + \gamma_t s_{i,t}$ where $W_t = \sum_{i \in V} w_{i,t}$ ,
$$\gamma_t = \frac{(1 + \beta) \eta}{\min_{i \in V} \sum_{j : j \xrightarrow{t} i} s_{j,t}} \qquad \text{and} \qquad \beta = 2\eta \sqrt{\frac{\ln(5K/\delta)}{\ln K}} \; ;$$

4. Play action $I_t$ drawn according to distribution $p_t = (p_{1,t}, \ldots, p_{K,t})$ ;

5. Observe pairs $(i, \ell_{i,t})$ for all $i \in S_{I_t, t}$;

6. For any $i \in V$ set $g_{i,t} = 1 - \ell_{i,t}$ and $w_{i,t+1} = w_{i,t} \exp(\eta \widehat{g}_{i,t})$, where
$$\widehat{g}_{i,t} = \frac{g_{i,t} \mathbb{I}\{i \in S_{I_t, t}\} + \beta}{q_{i,t}} \qquad \text{and} \qquad q_{i,t} = \sum_{j : j \xrightarrow{t} i} p_{j,t} \; .$$

---

**Theorem 9** *Let algorithm ELP.P run with learning rate $\eta \leq 1/(3K)$ sufficiently small such that $\beta \leq 1/4$. Then, with probability at least $1 - \delta$ we have*

$$L_{A,T} - \max_{k \in V} L_{k,T} \; \leq \; \sqrt{5 \ln\left(\frac{5}{\delta}\right) \sum_{t=1}^{T} \mathtt{mas}(G_t) + \frac{2 \ln(5K/\delta)}{\eta} + 12\eta \sqrt{\frac{\ln(5K/\delta)}{\ln K}} \sum_{t=1}^{T} \mathtt{mas}(G_t)}$$

$$+ \; \widetilde{\mathcal{O}}\left(1 + \sqrt{T}\eta + T\eta^2\right) \left(\max_{t=1 \ldots T} \mathtt{mas}^2(G_t)\right),$$

*where the $\widetilde{\mathcal{O}}$ notation hides only numerical constants and factors logarithmic in $K$ and $1/\delta$. In particular, if for constants $m_1, \ldots, m_T$ we have $\mathtt{mas}(G_t) \leq m_t$, $t = 1, \ldots, T$, and we pick $\eta$ such that*

$$\eta^2 = \frac{1}{6} \frac{\sqrt{\ln(5K/\delta)(\ln K)}}{\sum_{t=1}^{T} m_t}$$

*then we get the bound*

$$L_{A,T} - \max_{k \in V} L_{k,T} \leq 10 \frac{\ln^{1/4}(5K/\delta)}{\ln^{1/4} K} \sqrt{\ln\left(\frac{5K}{\delta}\right) \sum_{t=1}^{T} m_t} + \widetilde{\mathcal{O}}(T^{1/4}) \left(\max_{t=1...T} \mathtt{mas}^2(G_t)\right) \ .$$

This theorem essentially tells us that the regret of the ELP.P algorithm, up to second-order factors, is quantified by $\sqrt{\sum_{t=1}^{T} \mathtt{mas}(G_t)}$. Recall that, in the special case when $G_t$ is symmetric, we have $\mathtt{mas}(G_t) = \alpha(G_t)$.

One computational issue to bear in mind is that this theorem (as well as Theorem 2 and Corollary 3) holds under an optimal choice of $\eta$. In turn, this value depends on upper bounds on $\sum_{t=1}^{T} \mathtt{mas}(G_t)$ (or on $\sum_{t=1}^{T} \alpha(G_t)$, in the symmetric case). Unfortunately, in the worst case, computing the maximal acyclic subgraph or the independence number of a given graph is NP-hard, so implementing such algorithms is not *always* computationally tractable.[8] However, it is easy to see that the algorithm is robust to approximate computation of this value —misspecifying the average independence number $\frac{1}{T} \sum_{t=1}^{T} \alpha(G_t)$ by a factor of $v$ entails an additional $\sqrt{v}$ factor in the bound. Thus, one might use standard heuristics resulting in a reasonable approximation of the independence number. Although computing the independence number is also NP-hard to approximate, it is unlikely for intricate graphs with hard-to-approximate independence numbers to appear in relevant applications. Moreover, by setting the approximation to be either $K$ or 1, we trivially obtain an approximation factor of at most either $K$ or $\frac{1}{T} \sum_{t=1}^{T} \alpha(G_t)$. The former leads to a $\widetilde{\mathcal{O}}(\sqrt{KT})$ regret bound similar to the standard bandits setting, while the latter leads to a $\widetilde{\mathcal{O}}\left(\frac{1}{T} \sum_{t=1}^{T} \alpha(G_T)\sqrt{T}\right)$ regret bound, which is better than the regret for the bandits setting if the average independence number is less than $\sqrt{K}$. In contrast, this computational issue does not show up in Exp3-DOM, whose tuning relies only on efficiently-computable quantities.

## 5    Conclusions and Open Questions

In this paper we investigated online prediction problems in partial information regimes that interpolate between the classical bandit and expert settings. We provided algorithms, as well as upper and lower bounds on the attainable regret, with a non-trivial dependence on the information feedback structure. In particular, we have shown a number of results characterizing prediction performance in terms of: the structure of the feedback system, the amount of information available before prediction, and the nature (adversarial or fully random) of the process generating the feedback system.

There are many open questions that warrant further study, some of which are briefly mentioned below:

1. It would be interesting to study adaptations of our results to the case when the feedback system $\{S_{i,t}\}_{i \in V}$ may depend on the loss $\ell_{I_t,t}$ of player's action $I_t$. Note that this would prevent a direct construction of an unbiased estimator for unobserved losses, which many worst-case bandit algorithms (including ours —see the appendix) hinge upon.

---

[8] [20] proposed a generic mechanism to circumvent this, but the justification has a flaw which is not clear how to fix.

2. The upper bound contained in Theorem 2, expressed in terms of $\mathtt{mas}(\cdot)$, is almost certainly suboptimal, even in the uninformed setting, and it would be nice to see if more adequate graph complexity measures can be used instead.

3. Our lower bound in Theorem 5 refers to a constant graph sequence. We would like to provide a more complete characterization applying to sequences of adversarially-generated graphs $G_1, G_2, \ldots, G_T$ in terms of sequences of their corresponding independence numbers $\alpha(G_1), \alpha(G_2), \ldots, \alpha(G_T)$ (or variants thereof), in both the uninformed and the informed settings. Moreover, the adversary strategy achieving our lower bound is computationally hard to implement in the worst case (the adversary needs to identify the largest independent set in a given graph). What is the achievable regret if the adversary is assumed to be computationally bounded?

4. The information feedback models we used are natural and simple. They assume that the action at a give time period only affects rewards and observations for that period. In some settings, the reward observation may be delayed. In such settings, the action taken at a given stage may affect what is observed in subsequent stages. We leave the issue of modelling and analyzing such setting to future work.

5. Finally, we would like to see what is the achievable performance in the special case of stochastic rewards, which are assumed to be drawn i.i.d. from some unknown distributions. This was recently considered in [7], with results depending on the graph clique structure. However, the tightness of these results remains to be ascertained.

**Acknowledgments**

# References

[1] N. Alon, N. Cesa-Bianchi, C. Gentile, and Y. Mansour. From bandits to experts: A tale of domination and independence. In *NIPS*, 2013.

[2] N. Alon and J. H. Spencer. *The probabilistic method*. John Wiley & Sons, 2004.

[3] Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, 2009.

[4] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.

[5] Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.

[6] Y. Caro. New results on the independence number. In *Tech. Report, Tel-Aviv University*, 1979.

[7] Stéphane Caron, Branislav Kveton, Marc Lelarge, and Smriti Bhagat. Leveraging side observations in stochastic bandits. In *UAI*, 2012.

[8] N. Cesa-Bianchi, Y. Freund, D. Haussler, D. P. Helmbold, R. E. Schapire, and M. K. Warmuth. How to use expert advice. *J. ACM*, 44(3):427–485, 1997.

[9] N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.

[10] Nicolò Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *J. Comput. Syst. Sci.*, 78(5):1404–1422, 2012.

[11] Nicolò Cesa-Bianchi, Yishay Mansour, and Gilles Stoltz. Improved second-order bounds for prediction with expert advice. In *Proceedings of the 18th Annual Conference on Learning Theory*, pages 217–232, 2005.

[12] V. Chvatal. A greedy heuristic for the set-covering problem. *Mathematics of Operations Research*, 4(3):233–235, 1979.

[13] Ofer Dekel, Ambuj Tewari, and Raman Arora. Online bandit learning against an adaptive adversary: from regret to policy regret. In *ICML*, 2012.

[14] D.A. Freedman. On tail probabilities for martingales. *Annals of Probability*, 3:100–118, 1975.

[15] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In *Euro-COLT*, pages 23–37. Springer-Verlag, 1995. Also, JCSS 55(1): 119-139 (1997).

[16] A. M. Frieze. On the independence number of random graphs. *Discrete Mathematics*, 81:171–175, 1990.

[17] A. Kalai and S. Vempala. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences*, 71:291–307, 2005.

[18] T. Kocàk, G. Neu, M. Valko, and R. Munos. Efficient learning by implicit exploration in bandit problems with side observations. Manuscript, 2014.

[19] Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108:212–261, 1994.

[20] S. Mannor and O. Shamir. From bandits to experts: On the value of side-observations. In *25th Annual Conference on Neural Information Processing Systems (NIPS 2011)*, 2011.

[21] Alan Said, Ernesto W De Luca, and Sahin Albayrak. How social relationships affect user similarities. In *Proceedings of the International Conference on Intelligent User Interfaces Workshop on Social Recommender Systems, Hong Kong*, 2010.

[22] Gilles Stoltz. *Information Incomplète et Regret Interne en Prédiction de Suites Individuelles.* PhD thesis, Université Paris-XI Orsay, 2005.

[23] V. G. Vovk. Aggregating strategies. In *COLT*, pages 371–386, 1990.

[24] V. K. Wey. A lower bound on the stability number of a simple graph. In *Bell Lab. Tech. Memo No. 81-11217-9*, 1981.

## A    Technical lemmas and proofs from Section 3

This section contains the proofs of all technical results occurring in Section 3, along with ancillary graph-theoretic lemmas. Throughout this appendix, $\mathbb{E}_t[\cdot]$ is a shorthand for $\mathbb{E}\big[\cdot \mid I_1, \ldots, I_{t-1}\big]$. Also, for ease of exposition, we implicitly first condition on the history, i.e., $I_1, I_2, \ldots, I_{t-1}$, and later take an expectation with respect to that history. This implies that, given that conditioning, we can treat random variables such as $p_{i,t}$ as constants, and we can later take an expectation over history so as to remove the conditioning.

### A.1    Proof of Fact 4

Using standard properties of geometric sums, one can immediately see that

$$
\sum_{i=1}^{K} \frac{p_i}{\sum_{j=i}^{K} p_j} = \sum_{i=1}^{K-1} \frac{2^{-i}}{2^{-i+1}} + \frac{2^{-K+1}}{2^{-K+1}} = \frac{K-1}{2} + 1 = \frac{K+1}{2} \ ,
$$

hence the claimed result.

### A.2    Proof of Lemma 1

Following the proof of Exp3 [4], we have

$$
\begin{aligned}
\frac{W_{t+1}}{W_t} &= \sum_{i \in V} \frac{w_{i,t+1}}{W_t} \\
&= \sum_{i \in V} \frac{w_{i,t} \, \exp(-\eta \, \widehat{\ell}_{i,t})}{W_t} \\
&= \sum_{i \in V} p_{i,t} \, \exp(-\eta \, \widehat{\ell}_{i,t}) \\
&\leq \sum_{i \in V} p_{i,t} \left(1 - \eta \widehat{\ell}_{i,t} + \frac{1}{2} \eta^2 (\widehat{\ell}_{i,t})^2\right) \quad \text{using } e^{-x} \leq 1 - x + x^2/2 \text{ for all } x \geq 0 \\
&\leq 1 - \eta \sum_{i \in V} p_{i,t} \widehat{\ell}_{i,t} + \frac{\eta^2}{2} \sum_{i \in V} p_{i,t} (\widehat{\ell}_{i,t})^2 \ .
\end{aligned}
$$

Taking logs, using $\ln(1-x) \leq -x$ for all $x \geq 0$, and summing over $t = 1, \ldots, T$ yields

$$
\ln \frac{W_{T+1}}{W_1} \leq -\eta \sum_{t=1}^{T} \sum_{i \in V} p_{i,t} \widehat{\ell}_{i,t} + \frac{\eta^2}{2} \sum_{t=1}^{T} \sum_{i \in V} p_{i,t} (\widehat{\ell}_{i,t})^2 \ .
$$

Moreover, for any fixed comparison action $k$, we also have

$$\ln \frac{W_{T+1}}{W_1} \geq \ln \frac{w_{k,T+1}}{W_1} = -\eta \sum_{t=1}^{T} \widehat{\ell}_{k,t} - \ln K \ .$$

Putting together and rearranging gives

$$\sum_{t=1}^{T} \sum_{i \in V} p_{i,t} \widehat{\ell}_{i,t} \leq \sum_{t=1}^{T} \widehat{\ell}_{k,t} + \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \sum_{i \in V} p_{i,t} (\widehat{\ell}_{i,t})^2 \ . \tag{5}$$

Note that, for all $i \in V$,

$$\mathbb{E}_t[\widehat{\ell}_{i,t}] = \sum_{j : i \in S_{j,t}} p_{j,t} \frac{\ell_{i,t}}{q_{i,t}} = \sum_{j : j \xrightarrow{t} i} p_{j,t} \frac{\ell_{i,t}}{q_{i,t}} = \frac{\ell_{i,t}}{q_{i,t}} \sum_{j : j \xrightarrow{t} i} p_{j,t} = \ell_{i,t} \ .$$

Moreover,

$$\mathbb{E}_t\big[(\widehat{\ell}_{i,t})^2\big] = \sum_{j : i \in S_{j,t}} p_{j,t} \frac{\ell_{i,t}^2}{q_{i,t}^2} = \frac{\ell_{i,t}^2}{q_{i,t}^2} \sum_{j : j \xrightarrow{t} i} p_{j,t} \leq \frac{1}{q_{i,t}^2} \sum_{j : j \xrightarrow{t} i} p_{j,t} = \frac{1}{q_{i,t}} \ .$$

Hence, taking expectations $\mathbb{E}_t$ on both sides of (5), and recalling the definition of $Q_t$, we can write

$$\sum_{t=1}^{T} \sum_{i \in V} p_{i,t} \, \ell_{i,t} \leq \sum_{t=1}^{T} \ell_{k,t} + \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^{T} Q_t \ . \tag{6}$$

Finally, taking expectations over history to remove conditioning gives

$$\mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \mathbb{E}[Q_t]$$

as claimed. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad$ $\square$

### A.3   Proof of Theorem 2

We first need the following lemma.

**Lemma 10** *Let $G = (V, D)$ be a directed graph with vertex set $V = \{1, \ldots, K\}$, and arc set $D$. Then, for any distribution $p$ over $V$ we have,*

$$\sum_{i=1}^{K} \frac{p_i}{p_i + \sum_{j : j \to i} p_j} \leq \mathtt{mas}(G) \ .$$

**Proof.** We show that there is a subset of vertices $V'$ such that the induced graph is acyclic and $|V'| \geq \sum_{i=1}^{K} \frac{p_i}{p_i + \sum_{j \in N_i^-} p_j}$. Let $N_i^-$ be the in-neighborhood of node $i$, i.e., the set of nodes $j$ such that $(j, i) \in D$.

We prove the lemma by adding elements to an initially empty set $V'$. Let

$$\Phi_0 = \sum_{i=1}^{K} \frac{p_i}{p_i + \sum_{j\,:\,j\to i} p_j},$$

and let $i_1$ be the vertex which minimizes $p_i + \sum_{j\in N_i^-} p_j$ over $i \in V$. We now delete $i_1$ from the graph, along with all its incoming neighbors (set $N_{i_1}^-$), and all edges which are incident (both departing and incoming) to these nodes, and then iterating on the remaining graph. Let $N_{i,1}^-$ be the in-neighborhoods of the graph after the first step. The contribution of all the deleted vertices to $\Phi_0$ is

$$\sum_{r\in N_{i_1}^-\cup\{i_1\}} \frac{p_r}{p_r + \sum_{j\in N_r^-} p_j} \leq \sum_{r\in N_{i_1}^-\cup\{i_1\}} \frac{p_r}{p_{i_1} + \sum_{j\in N_{i_1}^-} p_j} = 1\,,$$

where the inequality follows from the minimality of $i_1$.

Let $V' \leftarrow V' \cup \{i_1\}$, and $V_1 = V \setminus (N_{i_1}^- \cup \{i_1\})$. Then, from the first step we have

$$\Phi_1 = \sum_{i\in V_1} \frac{p_i}{p_i + \sum_{j\in N_{i,1}^-} p_j} \geq \sum_{i\in V_1} \frac{p_i}{p_i + \sum_{j\in N_i^-} p_j} \geq \Phi_0 - 1\,.$$

We apply the very same argument to $\Phi_1$ with node $i_2$ (minimizing $p_i + \sum_{j\in N_{i,1}^-} p_j$ over $i \in V_1$), to $\Phi_2$ with node $i_3$, ..., to $\Phi_{s-1}$ with node $i_s$, up until $\Phi_s = 0$, i.e., until no nodes are left in the reduced graph. This gives $\Phi_0 \leq s = |V'|$, where $V' = \{i_1, i_2, \ldots, i_s\}$. Moreover, since in each step $r = 1, \ldots, s$ we remove all remaining arcs incoming to $i_r$, the graph induced by set $V'$ cannot contain cycles. $\qquad\square$

The claim of Theorem 2 follows from a direct combination of Lemma 1 with Lemma 10.

## A.4 Proof of Theorem 5

The proof uses a variant of the standard multi-armed bandit lower bound [9]. The intuition is that when we have $\alpha(G)$ non-adjacent nodes, the problem reduces to an instance of the standard multi-armed bandit (where information beyond the loss of the action choses is observed) on $\alpha(G)$ actions.

By Yao's minimax principle, in order to establish the lower bound, it is enough to demonstrate some probabilistic adversary strategy, on which the expected regret of any *deterministic* algorithm $A$ is bounded from below by $C\sqrt{\alpha(G)T}$ for some constant $C$.

Specifically, suppose without loss of generality that we number the nfiodes in some largest independent set of $G$ by $1, 2, \ldots, \alpha(G)$, and all the other nodes in the graph by $\alpha(G) + 1, \ldots, |V|$. Let $\epsilon$ be a parameter to be determined later, and consider the following joint distribution over stochastic loss sequences:

- Let $Z$ be uniformly distributed on $1, 2, \ldots, \alpha(G)$;

- Conditioned on $Z = i$, each loss $\ell_{j,t}$ is independent Bernoulli with parameter $1/2$ if $j \neq i$ and $j < \alpha(G)$, independent Bernoulli with parameter $1/2 - \epsilon$ if $j = i$, and is 1 with probability 1, otherwise.

For each $i = 1 \ldots \alpha(G)$, let $T_i$ be the number of times the node $i$ was chosen by the algorithm after $T$ rounds. Also, let $T_\Delta$ denote the number of times some node whose index is larger than $\alpha(G)$ is chosen after $T$ rounds. Finally, let $\mathbb{E}_i$ denote expectation conditioned on $Z = i$, and $\mathbb{P}_i$ denote the probability over loss sequences conditioned on $Z = i$. We have

$$
\max_{k \in V} \mathbb{E}[L_{A,T} - L_{k,T}] = \frac{1}{\alpha(G)} \sum_{i=1}^{\alpha(G)} \mathbb{E}_i \left[ L_{A,T} - \left( \frac{1}{2} - \epsilon \right) T \right]
$$

$$
= \frac{1}{\alpha(G)} \sum_{i=1}^{\alpha(G)} \mathbb{E}_i \left[ \sum_{j \in \{1 \ldots \alpha(G)\} \setminus i} \frac{1}{2} T_j + \left( \frac{1}{2} - \epsilon \right) T_i + T_\Delta - \left( \frac{1}{2} - \epsilon \right) T \right]
$$

$$
= \frac{1}{\alpha(G)} \sum_{i=1}^{\alpha(G)} \mathbb{E}_i \left[ \frac{1}{2} \sum_{j=1}^{\alpha(G)} T_j + \frac{1}{2} T_\Delta + \frac{1}{2} T_\Delta - \epsilon T_i - \left( \frac{1}{2} - \epsilon \right) T \right].
$$

Since $\sum_{j=1}^{\alpha(G)} T_j + T_\Delta = T$, this expression equals

$$
\frac{1}{\alpha(G)} \sum_{i=1}^{\alpha(G)} \mathbb{E}_i \left[ \frac{1}{2} T_\Delta + \epsilon(T - T_i) \right] \;\geq\; \epsilon \left( T - \frac{1}{\alpha(G)} \sum_{i=1}^{\alpha(G)} \mathbb{E}_i[T_i] \right). \tag{7}
$$

Now, consider another distribution $\mathbb{P}_0$ over the loss sequence, which corresponds to the distribution above but with $\epsilon = 0$ (namely, all nodes $1, \ldots, \alpha(G)$ have losses which are $\pm 1$ independently and with equal probability, and all nodes whose index is larger than $\alpha(G)$ have losses of 1), and denote by $\mathbb{E}_0$ the corresponding expectation. We upper bound the difference between $\mathbb{E}_i[T_i]$ and $\mathbb{E}_0[T_i]$, using information theoretic arguments. Let $\lambda_t$ be the collection of loss values observed at round $t$, and $\lambda^t = (\lambda_1, \ldots, \lambda_t)$. Note that since the algorithm is deterministic, $\lambda^{t-1}$ determines the algorithm's choice of action $I_t$ at each round $t$, and hence $T_i$ is determined by $\lambda^T$, and thus $\mathbb{E}_0[T_i \mid \lambda^T] = \mathbb{E}_i[T_i \mid \lambda^T]$. We have

$$
\mathbb{E}_i[T_i] - \mathbb{E}_0[T_i] = \sum_{\lambda^T} \mathbb{P}_i(\lambda^T) \mathbb{E}_i[T_i \mid \lambda^T] - \sum_{\lambda^T} \mathbb{P}_0(\lambda^T) \mathbb{E}_0[T_i \mid \lambda^T]
$$

$$
= \sum_{\lambda^T} \mathbb{P}_i(\lambda^T) \mathbb{E}_i[T_i \mid \lambda^T] - \sum_{\lambda^T} \mathbb{P}_0(\lambda^T) \mathbb{E}_i[T_i \mid \lambda^T]
$$

$$
\leq \sum_{\lambda^T : \mathbb{P}_i(\lambda^T) > \mathbb{P}_0(\lambda^T)} \left( \mathbb{P}_i(\lambda^T) - \mathbb{P}_0(\lambda^T) \right) \mathbb{E}_i[T_i \mid \lambda^T]
$$

$$
\leq T \sum_{\lambda^T : \mathbb{P}_i(\lambda^T) > \mathbb{P}_0(\lambda^T)} \left( \mathbb{P}_i(\lambda^T) - \mathbb{P}_0(\lambda^T) \right).
$$

Using Pinsker's inequality, this is at most

$$
T \sqrt{\frac{1}{2} D_{\mathrm{kl}}(\mathbb{P}_0(\lambda^T) \,\|\, \mathbb{P}_i(\lambda^T))}
$$

where $D_{\mathrm{kl}}$ is the Kullback-Leibler divergence (or relative entropy) between the distributions $\mathbb{P}_i$ and $\mathbb{P}_0$. Using the chain rule for relative entropy, this equals

$$
T \sqrt{\frac{1}{2} \sum_{t=1}^{T} \sum_{\lambda^{t-1}} \mathbb{P}_0(\lambda^{t-1}) D_{\mathrm{kl}}\left( \mathbb{P}_0(\lambda_t | \lambda^{t-1}) \,\|\, \mathbb{P}_i(\lambda_t | \lambda^{t-1}) \right)}.
$$

Let us consider any single relative entropy term above. Recall that $\lambda^{t-1}$ determines the node $I_t$ picked at round $t$. If this node is not $i$ or adjacent to $i$, then $\lambda_t$ is going to have the same distribution under both $\mathbb{P}_i$ and $\mathbb{P}_0$, and the relative entropy is zero. Otherwise, the coordinate of $\lambda_t$ corresponding to node $i$ (and that coordinate only) will have a different distribution: Bernoulli with parameter $\frac{1}{2} - \epsilon$ under $\mathbb{P}_i$, and Bernoulli with parameter $\frac{1}{2}$ under $\mathbb{P}_0$. The relative entropy term in this case is easily shown to be $-\frac{1}{2}\log(1 - 4\epsilon^2) \le 8\log(4/3)\,\epsilon^2$. Therefore, letting $S_{I_t}$ denote the feedback set at time $t$, we can upper bound the above by

$$T\sqrt{\frac{1}{2}\sum_{t=1}^{T}\mathbb{P}_0(i \in S_{I_t})(8\log(4/3)\epsilon^2)} \;=\; 2T\epsilon\sqrt{\log\left(\frac{4}{3}\right)\mathbb{E}_0\big[|\{t : i \in S_{I_t}\}|\big]}$$

$$\le\; 2T\epsilon\sqrt{\log\left(\frac{4}{3}\right)\mathbb{E}_0\left[T_i + T_\Delta\right]}. \qquad (8)$$

We now claim that we can assume $\mathbb{E}_0[T_\Delta] \le 0.08\sqrt{\alpha(G)T}$. To see why, note that if $\mathbb{E}_0[T_\Delta] > 0.08\sqrt{\alpha(G)T}$, then the expected regret under $\mathbb{E}_0$ would have been at least

$$\max_{k \in V}\mathbb{E}_0[L_{A,T} - L_{k,T}] \;=\; \mathbb{E}_0\left[T_\Delta + \frac{1}{2}\sum_{j=1}^{\alpha(G)}T_j\right] - \frac{1}{2}T$$

$$=\; \mathbb{E}_0\left[\frac{1}{2}T_\Delta + \frac{1}{2}\left(T_\Delta + \sum_{j=1}^{\alpha(G)}T_j\right)\right] - \frac{1}{2}T$$

$$=\; \mathbb{E}_0\left[\frac{1}{2}T_\Delta + \frac{1}{2}T\right] - \frac{1}{2}T$$

$$=\; \frac{1}{2}\mathbb{E}_0[T_\Delta]$$

$$>\; 0.04\sqrt{\alpha(G)T}.$$

So for the adversary strategy defined by the distribution $\mathbb{P}_0$, we would get an expected regret lower bound as required. Thus, it only remains to treat the case where $\mathbb{E}_0[T_\Delta] \le 0.08\sqrt{\alpha(G)T}$. Plugging in this upper bound into Eq. (8), we get overall that

$$\mathbb{E}_i[T_i] - \mathbb{E}_0[T_i] \le 2T\epsilon\sqrt{\log\left(\frac{4}{3}\right)\mathbb{E}_0\left[T_i + 0.08\sqrt{\alpha(G)T}\right]}.$$

Therefore, the expected regret lower bound in Eq. (7) is at least

$$\epsilon\left(T - \frac{1}{\alpha(G)}\sum_{i=1}^{\alpha(G)}\mathbb{E}_0[T_i] - \frac{1}{\alpha(G)}\sum_{i=1}^{\alpha(G)}2T\epsilon\sqrt{\log\left(\frac{4}{3}\right)\mathbb{E}_0\left[T_i + 0.08\sqrt{\alpha(G)T}\right]}\right)$$

$$\ge\; \epsilon\left(T - \frac{T}{\alpha(G)} - 2T\epsilon\sqrt{\log\left(\frac{4}{3}\right)\frac{1}{\alpha(G)}\sum_{i=1}^{\alpha(G)}\mathbb{E}_0\left[T_i + 0.08\sqrt{\alpha(G)T}\right]}\right)$$

$$\ge\; \epsilon T\left(1 - \frac{1}{\alpha(G)} - 2\epsilon\sqrt{\log\left(\frac{4}{3}\right)\left(\frac{T}{\alpha(G)} + 0.08\sqrt{\alpha(G)T}\right)}\right).$$

Since $\alpha(G) > 1$, we have $1 - \frac{1}{\alpha(G)} \geq \frac{1}{2}$, and since $T \geq 0.0064\alpha^3(G)$, we have $0.08\sqrt{\alpha(G)T} \leq \frac{T}{\alpha(G)}$. Overall, we can lower bound the expression above by

$$\epsilon T \left( \frac{1}{2} - 2\epsilon \sqrt{2\log\left(\frac{4}{3}\right) \frac{T}{\alpha(G)}} \right) .$$

Picking $\epsilon = \frac{1}{8\sqrt{2\log(4/3)T/\alpha(G)}}$, the expression above is

$$\frac{T}{8\sqrt{2\log\left(\frac{4}{3}\right) \frac{T}{\alpha(G)}}} \frac{1}{4} \geq 0.04\sqrt{\alpha(G)T} .$$

This constitutes a lower bound on the expected regret, from which the result follows.

## A.5   Proof of Theorem 6

Fix round $t$, and let $G = (V, D)$ be the Erdős-Renyi random graph generated at time $t$, $N_i^-$ be the in-neighborhood of node $i$, i.e., the set of nodes $j$ such that $(j, i) \in D$, and denote by $d_i^-$ the indegree of $i$. We need the following lemmas.

**Lemma 11** *Fix a directed graph $G = (V, D)$. Let $p_1, \ldots, p_K$ be an arbitrary probability distribution defined over $V$, $f : V \to V$ be an arbitrary permutation of $V$, and $\mathbb{E}_f$ denote the expectation w.r.t. a random permutation $f$. Then, for any $i \in V$, we have*

$$\mathbb{E}_f \left[ \frac{p_{f(i)}}{p_{f(i)} + \sum_{j : j \to i} p_{f(j)}} \right] = \frac{1}{1 + d_i^-} .$$

**Proof.**  Consider selecting a subset $S \subset V$ of $1 + d_i^-$ nodes. We consider the contribution to the expectation when $S = N_{f(i)}^- \cup \{f(i)\}$. Since there are $K(K-1) \cdots (K - d_i^- + 1)$ terms (out of $K!$) contributing to the expectation, we can write

$$\mathbb{E}_f \left[ \frac{p_{f(i)}}{p_{f(i)} + \sum_{j : f(j) \in N_{f(i)}^-} p_{f(j)}} \right] = \frac{1}{\binom{K}{d_i^-}} \sum_{S \subset V, |S| = d_i^-} \frac{1}{1 + d_i^-} \sum_{i \in S} \frac{p_i}{p_i + \sum_{j \in S, j \neq i} p_j}$$

$$= \frac{1}{\binom{K}{d_i^-}} \sum_{S \subset V, |S| = d_i^-} \frac{1}{1 + d_i^-}$$

$$= \frac{1}{1 + d_i^-} .$$

$\square$

**Lemma 12** *Let $p_1, \ldots, p_K$ be an arbitrary probability distribution defined over $V$, and $\mathbb{E}$ denote the expectation w.r.t. the Erdős-Renyi random draw of arcs at time $t$. Then, for any fixed $i \in V$, we have*

$$\mathbb{E} \left[ \frac{p_i}{p_i + \sum_{j : j \xrightarrow{t} i} p_j} \right] = \frac{1}{rK} \left( 1 - (1 - r)^K \right) .$$

**Proof.** For the given $i \in V$ and time $t$, consider the Bernoulli random variables $X_j, j \in V \setminus \{i\}$, and denote by $\mathbb{E}_{j\,:\,j\neq i}$ the expectation w.r.t. all of them. We symmetrize $\mathbb{E}\left[\frac{p_i}{p_i + \sum_{j\,:\,j\xrightarrow{t}i} p_j}\right]$ by means of a random permutation $f$, as in Lemma 11. We can write

$$
\begin{aligned}
\mathbb{E}\left[\frac{p_i}{p_i + \sum_{j\,:\,j\xrightarrow{t}i} p_j}\right] &= \mathbb{E}_{j\,:\,j\neq i}\left[\frac{p_i}{p_i + \sum_{j\,:\,j\neq i} X_j p_j}\right] \\
&= \mathbb{E}_{j\,:\,j\neq i}\mathbb{E}_f\left[\frac{p_{f(i)}}{p_{f(i)} + \sum_{j\,:\,j\neq i} X_{f(j)} p_{f(j)}}\right] \qquad \text{(by symmetry)} \\
&= \mathbb{E}_{j\,:\,j\neq i}\left[\frac{1}{1 + \sum_{j\,:\,j\neq i} X_j}\right] \qquad \text{(from Lemma 11)} \\
&= \sum_{i=0}^{K-1}\binom{K-1}{i} r^i (1-r)^{K-1-i}\frac{1}{i+1} \\
&= \frac{1}{rK}\sum_{i=0}^{K-1}\binom{K}{i+1} r^{i+1}(1-r)^{K-1-i} \\
&= \frac{1}{rK}\left(1 - (1-r)^K\right) .
\end{aligned}
$$

$\square$

At this point, we follow the proof of Lemma 1 up until (6). We take an expectation $\mathbb{E}_{G_1,\dots,G_T}$ w.r.t. the randomness in generating the sequence of graphs $G_1,\dots,G_T$. This yields

$$
\sum_{t=1}^{T}\mathbb{E}_{G_1,\dots,G_T}\left[\sum_{i\in V} p_{i,t}\,\ell_{i,t}\right] \leq \sum_{t=1}^{T}\ell_{k,t} + \frac{\ln K}{\eta} + \frac{\eta}{2}\sum_{t=1}^{T}\mathbb{E}_{G_1,\dots,G_T}\left[Q_t\right] .
$$

We use Lemma 12 to upper bound $\mathbb{E}_{G_1,\dots,G_T}\left[Q_t\right]$ by $\frac{1}{r}\left(1 - (1-r)^K\right)$, and take the outer expectation to remove conditioning, as in the proof of Lemma 1. This concludes the proof.

# B  Technical lemmas and proofs from Section 4.1

Again, throughout this appendix, $\mathbb{E}_t[\cdot]$ is a shorthand for the conditional expectation $\mathbb{E}_t[\cdot \mid I_1, I_2, \dots, I_{t-1}]$. Moreover, as we did in Appendix A, in round $t$ we first condition on the history $I_1, I_2, \dots, I_{t-1}$, and then take an outer expectation with respect to that history.

## B.1  Proof of Lemma 7

We start to bound the contribution to the overall regret of an instance indexed by $b$. When clear from the context, we remove the superscript $b$ from $\gamma^{(b)}$, $w_{i,t}^{(b)}$, $p_{i,t}^{(b)}$, and other related quantities.

For any $t \in T^{(b)}$ we have

$$
\begin{aligned}
\frac{W_{t+1}}{W_t} &= \sum_{i \in V} \frac{w_{i,t+1}}{W_t} \\
&= \sum_{i \in V} \frac{w_{i,t}}{W_t} \exp\left(-(\gamma/2^b)\,\widehat{\ell}_{i,t}\right) \\
&= \sum_{i \in R_t} \frac{p_{i,t} - \gamma/|R_t|}{1-\gamma} \exp\left(-(\gamma/2^b)\,\widehat{\ell}_{i,t}\right) + \sum_{i \notin R_t} \frac{p_{i,t}}{1-\gamma} \exp\left(-(\gamma/2^b)\,\widehat{\ell}_{i,t}\right) \\
&\leq \sum_{i \in R_t} \frac{p_{i,t} - \gamma/|R_t|}{1-\gamma} \left(1 - \frac{\gamma}{2^b}\widehat{\ell}_{i,t} + \frac{1}{2}\left(\frac{\gamma}{2^b}\widehat{\ell}_{i,t}\right)^2\right) + \sum_{i \notin R_t} \frac{p_{i,t}}{1-\gamma} \left(1 - \frac{\gamma}{2^b}\widehat{\ell}_{i,t} + \frac{1}{2}\left(\frac{\gamma}{2^b}\widehat{\ell}_{i,t}\right)^2\right)
\end{aligned}
$$

(using $e^{-x} \leq 1 - x + x^2/2$ for all $x \geq 0$)

$$
\leq 1 - \frac{\gamma/2^b}{1-\gamma}\sum_{i \in V}p_{i,t}\widehat{\ell}_{i,t} + \frac{\gamma^2/2^b}{1-\gamma}\sum_{i \in R_t}\frac{\widehat{\ell}_{i,t}}{|R_t|} + \frac{1}{2}\frac{(\gamma/2^b)^2}{1-\gamma}\sum_{i \in V}p_{i,t}\left(\widehat{\ell}_{i,t}\right)^2 .
$$

Taking logs, upper bounding, and summing over $t \in T^{(b)}$ yields

$$
\ln \frac{W_{|T^{(b)}|+1}}{W_1} \leq -\frac{\gamma/2^b}{1-\gamma}\sum_{t \in T^{(b)}}\sum_{i \in V}p_{i,t}\widehat{\ell}_{i,t} + \frac{\gamma^2/2^b}{1-\gamma}\sum_{t \in T^{(b)}}\sum_{i \in R_t}\frac{\widehat{\ell}_{i,t}}{|R_t|} + \frac{1}{2}\frac{(\gamma/2^b)^2}{1-\gamma}\sum_{t \in T^{(b)}}\sum_{i \in V}p_{i,t}\left(\widehat{\ell}_{i,t}\right)^2 .
$$

Moreover, for any fixed comparison action $k$, we also have

$$
\ln \frac{W_{|T^{(b)}|+1}}{W_1} \geq \ln \frac{w_{k,|T^{(b)}|+1}}{W_1} = -\frac{\gamma}{2^b}\sum_{t \in T^{(b)}}\widehat{\ell}_{k,t} - \ln K .
$$

Putting together, rearranging, and using $1 - \gamma \leq 1$ gives

$$
\sum_{t \in T^{(b)}}\sum_{i \in V}p_{i,t}\widehat{\ell}_{i,t} \leq \sum_{t \in T^{(b)}}\widehat{\ell}_{k,t} + \frac{2^b \ln K}{\gamma} + \gamma\sum_{t \in T^{(b)}}\sum_{i \in R_t}\frac{\widehat{\ell}_{i,t}}{|R_t|} + \frac{\gamma}{2^{b+1}}\sum_{t \in T^{(b)}}\sum_{i \in V}p_{i,t}\left(\widehat{\ell}_{i,t}\right)^2 .
$$

Reintroducing the notation $\gamma^{(b)}$ and summing over $b = 0, 1, \ldots, \lfloor \log_2 K \rfloor$ gives

$$
\sum_{t=1}^{T}\left(\sum_{i \in V}p_{i,t}^{(b_t)}\widehat{\ell}_{i,t}^{(b_t)} - \widehat{\ell}_{k,t}\right) \leq \sum_{b=0}^{\lfloor \log_2 K \rfloor}\frac{2^b \ln K}{\gamma^{(b)}} + \sum_{t=1}^{T}\sum_{i \in R_t}\frac{\gamma^{(b_t)}\widehat{\ell}_{i,t}^{(b_t)}}{|R_t|} + \sum_{t=1}^{T}\frac{\gamma^{(b_t)}}{2^{b_t+1}}\sum_{i \in V}p_{i,t}^{(b_t)}\left(\widehat{\ell}_{i,t}^{(b_t)}\right)^2 . \quad (9)
$$

Now, similarly to the proof of Lemma 1, we have that $\mathbb{E}_t\left[\widehat{\ell}_{i,t}^{(b_t)}\right] = \ell_{i,t}$ and $\mathbb{E}_t\left[(\widehat{\ell}_{i,t}^{(b_t)})^2\right] \leq \frac{1}{q_{i,t}^{(b_t)}}$ for any $i$ and $t$. Hence, taking expectations $\mathbb{E}_t$ on both sides of (9) and recalling the definition of $Q_t^{(b)}$ gives

$$
\sum_{t=1}^{T}\left(\sum_{i \in V}p_{i,t}^{(b_t)}\ell_{i,t} - \ell_{k,t}\right) \leq \sum_{b=0}^{\lfloor \log_2 K \rfloor}\frac{2^b \ln K}{\gamma^{(b)}} + \sum_{t=1}^{T}\sum_{i \in R_t}\frac{\gamma^{(b_t)}\ell_{i,t}}{|R_t|} + \sum_{t=1}^{T}\frac{\gamma^{(b_t)}}{2^{b_t+1}}Q_t^{(b_t)} . \quad (10)
$$

Moreover,

$$
\sum_{t=1}^{T}\sum_{i \in R_t}\frac{\gamma^{(b_t)}\ell_{i,t}}{|R_t|} \leq \sum_{t=1}^{T}\sum_{i \in R_t}\frac{\gamma^{(b_t)}}{|R_t|} = \sum_{t=1}^{T}\gamma^{(b_t)} = \sum_{b=0}^{\lfloor \log_2 K \rfloor}\gamma^{(b)}|T^{(b)}|
$$

and

$$\sum_{t=1}^{T} \frac{\gamma^{(b_t)}}{2^{b_t+1}} Q_t^{(b_t)} = \sum_{b=0}^{\lfloor \log_2 K \rfloor} \frac{\gamma^{(b)}}{2^{b+1}} \sum_{t \in T^{(b)}} Q_t^{(b)} \, .$$

Hence, substituting back into (10), taking outer expectations on both sides and recalling that $T^{(b)}$ is a random variable (since the adversary adaptively decides which steps $t$ fall into $T^{(b)}$), we get

$$\mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq \sum_{b=0}^{\lfloor \log_2 K \rfloor} \mathbb{E}\left[ \frac{2^b \ln K}{\gamma^{(b)}} + \gamma^{(b)} |T^{(b)}| + \frac{\gamma^{(b)}}{2^{b+1}} \sum_{t \in T^{(b)}} Q_t^{(b)} \right]$$

$$= \sum_{b=0}^{\lfloor \log_2 K \rfloor} \left( \frac{2^b \ln K}{\gamma^{(b)}} + \gamma^{(b)} \mathbb{E}\left[ \sum_{t \in T^{(b)}} \left( 1 + \frac{Q_t^{(b)}}{2^{b+1}} \right) \right] \right) \, . \tag{11}$$

This establishes (3).

In order to prove inequality (4), we need to tune each $\gamma^{(b)}$ separately. However, a good choice of $\gamma^{(b)}$ depends on the unknown random quantity

$$\overline{Q}^{(b)} = \sum_{t \in T^{(b)}} \left( 1 + \frac{Q_t^{(b)}}{2^{b+1}} \right) \, .$$

To overcome this problem, we slightly modify Exp3-DOM by applying a doubling trick[9] to guess $\overline{Q}^{(b)}$ for each $b$. Specifically, for each $b = 0, 1, \ldots, \lfloor \log_2 K \rfloor$, we use a sequence $\gamma_r^{(b)} = \sqrt{(2^b \ln K)/2^r}$, for $r = 0, 1, \ldots$. We initially run the algorithm with $\gamma_0^{(b)}$. Whenever the algorithm is running with $\gamma_r^{(b)}$ and observes that $\sum_s \overline{Q}_s^{(b)} > 2^r$, where the sum is over all $s$ so far in $T^{(b)}$,[10] then we restart the algorithm with $\gamma_{r+1}^{(b)}$. Because the contribution of instance $b$ to (11) is

$$\frac{2^b \ln K}{\gamma^{(b)}} + \gamma^{(b)} \sum_{t \in T^{(b)}} \left( 1 + \frac{Q_t^{(b)}}{2^{b+1}} \right)$$

the regret we pay when using any $\gamma_r^{(b)}$ is at most $2\sqrt{(2^b \ln K) 2^r}$. The largest $r$ we need is $\lceil \log_2 \overline{Q}^{(b)} \rceil$ and

$$\sum_{r=0}^{\lceil \log_2 \overline{Q}^{(b)} \rceil} 2^{r/2} < 5\sqrt{\overline{Q}^{(b)}} \, .$$

Since we pay regret at most 1 for each restart, we get

$$\mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq c \sum_{b=0}^{\lfloor \log_2 K \rfloor} \mathbb{E}\left[ \sqrt{(\ln K) \left( 2^b |T^{(b)}| + \frac{1}{2} \sum_{t \in T^{(b)}} Q_t^{(b)} \right)} + \lceil \log_2 \overline{Q}^{(b)} \rceil \right] \, .$$

---

[9] The pseudo-code for the variant of Exp3-DOM using such a doubling trick is not displayed here, since it is by now a folklore technique.

[10] Notice that $\sum_s \overline{Q}_s^{(b)}$ is an observable quantity.

for some positive constant $c$. Taking into account that

$$\sum_{b=0}^{\lfloor \log_2 K \rfloor} 2^b |T^{(b)}| \leq 2 \sum_{t=1}^{T} |R_t|$$

$$\sum_{b=0}^{\lfloor \log_2 K \rfloor} \sum_{t \in T^{(b)}} Q_t^{(b)} = \sum_{t=1}^{T} Q_t^{(b_t)}$$

$$\sum_{b=0}^{\lfloor \log_2 K \rfloor} \left\lceil \log_2 \overline{Q}^{(b)} \right\rceil = \mathcal{O}\big((\ln K) \ln(KT)\big)$$

we obtain

$$\mathbb{E}\big[L_{A,T} - L_{k,T}\big] \leq c \sum_{b=0}^{\lfloor \log_2 K \rfloor} \mathbb{E}\left[ \sqrt{(\ln K)\left(2^b |T^{(b)}| + \frac{1}{2}\sum_{t \in T^{(b)}} Q_t^{(b)}\right)} \right] + \mathcal{O}\big((\ln K) \ln(KT)\big)$$

$$\leq c \lfloor \log_2 K \rfloor \mathbb{E}\left[ \sqrt{\frac{\ln K}{\lfloor \log_2 K \rfloor} \sum_{t=1}^{T} \left(2|R_t| + \frac{1}{2}Q_t^{(b_t)}\right)} \right] + \mathcal{O}\big((\ln K) \ln(KT)\big)$$

$$= \mathcal{O}\left( (\ln K) \mathbb{E}\left[ \sqrt{\sum_{t=1}^{T} \left(4|R_t| + Q_t^{(b_t)}\right)} \right] + (\ln K) \ln(KT) \right)$$

as desired.

## B.2 Proof of Theorem 8

The following graph-theoretic lemma turns out to be fairly useful for analyzing directed settings. It is a directed-graph counterpart to a well-known result [6, 24] holding for undirected graphs.

**Lemma 13** *Let $G = (V, D)$ be a directed graph, with $V = \{1, \ldots, K\}$. Let $d_i^-$ be the indegree of node $i$, and $\alpha = \alpha(G)$ be the independence number of $G$. Then*

$$\sum_{i=1}^{K} \frac{1}{1 + d_i^-} \leq 2\alpha \ln\left(1 + \frac{K}{\alpha}\right) .$$

**Proof.** We proceed by induction, starting from the original $K$-node graph $G = G_K$ with indegrees $\{d_i^-\}_{i=1}^{K} = \{d_{i,K}^-\}_{i=1}^{K}$, and independence number $\alpha = \alpha_K$, and then progressively reduce $G$ by eliminating nodes and incident (both departing and incoming) arcs, thereby obtaining a sequence of smaller and smaller graphs $G_K, G_{K-1}, G_{K-2}, \ldots$, associated indegrees $\{d_{i,K}^-\}_{i=1}^{K}$, $\{d_{i,K-1}^-\}_{i=1}^{K-1}$, $\{d_{i,K-2}^-\}_{i=1}^{K-2}$, $\ldots$, and independence numbers $\alpha_K, \alpha_{K-1}, \alpha_{K-2}, \ldots$. Specifically, in step $s$ we sort nodes $i = 1, \ldots, s$ of $G_s$ in nonincreasing value of $d_{i,s}^-$, and obtain $G_{s-1}$ from $G_s$ by eliminating node 1 (i.e., the one having the largest indegree among the nodes of $G_s$), along with its incident arcs. On all such graphs, we use the classical Turan's theorem (e.g., [2]) stating that any *undirected*

27

graph with $n_s$ nodes and $m_s$ edges has an independent set of size at least $\frac{n_s}{\frac{2m_s}{n_s}+1}$. This implies that if $G_s = (V_s, D_s)$, then $\alpha_s$ satisfies[11]

$$\frac{|D_s|}{|V_s|} \geq \frac{|V_s|}{2\alpha_s} - \frac{1}{2} . \tag{12}$$

We then start from $G_K$. We can write

$$d^-_{1,K} = \max_{i=1\dots K} d^-_{i,K} \geq \frac{1}{K} \sum_{i=1}^{K} d^-_{i,K} = \frac{|D_K|}{|V_K|} \geq \frac{|V_K|}{2\alpha_K} - \frac{1}{2} .$$

Hence,

$$
\begin{aligned}
\sum_{i=1}^{K} \frac{1}{1+d^-_{i,K}} &= \frac{1}{1+d^-_{1,K}} + \sum_{i=2}^{K} \frac{1}{1+d^-_{i,K}} \\
&\leq \frac{2\alpha_K}{\alpha_K + K} + \sum_{i=2}^{K} \frac{1}{1+d^-_{i,K}} \\
&\leq \frac{2\alpha_K}{\alpha_K + K} + \sum_{i=1}^{K-1} \frac{1}{1+d^-_{i,K-1}}
\end{aligned}
$$

where the last inequality follows from $d^-_{i+1,K} \geq d^-_{i,K-1}$, $i = 1, \dots K-1$, due to the arc elimination trasforming $G_K$ into $G_{K-1}$. Recursively applying the same argument to $G_{K-1}$ (i.e., to the sum $\sum_{i=1}^{K-1} \frac{1}{1+d^-_{i,K-1}}$), and then iterating all the way to $G_1$ yields the upper bound

$$\sum_{i=1}^{K} \frac{1}{1+d^-_{i,K}} \leq \sum_{i=1}^{K} \frac{2\alpha_i}{\alpha_i + i} .$$

Combining with $\alpha_i \leq \alpha_K = \alpha$, and $\sum_{i=1}^{K} \frac{1}{\alpha+i} \leq \ln\left(1 + \frac{K}{\alpha}\right)$ concludes the proof. $\qquad\square$

The next lemma relates the size $|R_t|$ of the dominating set $R_t$ computed by the Greedy Set Cover algorithm of [12], operating on the time-$t$ feedback system $\{S_{i,t}\}_{i \in V}$, to the independence number $\alpha(G_t)$ and the domination number $\gamma(G_t)$ of $G_t$.

**Lemma 14** *Let $\{S_i\}_{i \in V}$ be a feedback system, and $G = (V, D)$ be the induced directed graph, with vertex set $V = \{1, \dots, K\}$, independence number $\alpha = \alpha(G)$, and domination number $\gamma = \gamma(G)$. Then the dominating set $R$ constructed by the Greedy Set Cover algorithm (see Section 2) satisfies*

$$|R| \leq \min\big\{\gamma(1 + \ln K), \lceil 2\alpha \ln K \rceil + 1\big\} .$$

**Proof.** As recalled in Section 2, the Greedy Set Cover algorithm of [12] achieves $|R| \leq \gamma(1+\ln K)$. In order to prove the other bound, consider the sequence of graphs $G = G_1, G_2, \dots$, where each $G_{s+1} = (V_{s+1}, D_{s+1})$ is obtained by removing from $G_s$ the vertex $i_s$ selected by the Greedy Set

---

[11] Note that $|D_s|$ is at least as large as the number of edges of the undirected version of $G_s$ which the independence number $\alpha_s$ actually refers to.

Cover algorithm, together with all the vertices in $G_s$ that are dominated by $i_s$, and all arcs incident to these vertices. By definition of the algorithm, the outdegree $d_s^+$ of $i_s$ in $G_s$ is largest in $G_s$. Hence,

$$d_s^+ \geq \frac{|D_s|}{|V_s|} \geq \frac{|V_s|}{2\alpha_s} - \frac{1}{2} \geq \frac{|V_s|}{2\alpha} - \frac{1}{2}$$

by Turan's theorem (e.g., [2]), where $\alpha_s$ is the independence number of $G_s$ and $\alpha \geq \alpha_s$. This shows that

$$|V_{s+1}| = |V_s| - d_s^+ - 1 \leq |V_s|\left(1 - \frac{1}{2\alpha}\right) \leq |V_s| e^{-1/(2\alpha)} .$$

Iterating, we obtain $|V_s| \leq K\, e^{-s/(2\alpha)}$. Choosing $s = \lceil 2\alpha \ln K \rceil + 1$ gives $|V_s| < 1$, thereby covering all nodes. Hence the dominating set $R = \{i_1, \ldots, i_s\}$ so constructed satisfies $|R| \leq \lceil 2\alpha \ln K \rceil + 1$.
$\square$

**Lemma 15** *If $a, b \geq 0$, and $a + b \geq B > A > 0$, then*

$$\frac{a}{a+b-A} \leq \frac{a}{a+b} + \frac{A}{B-A} .$$

**Proof.**

$$\frac{a}{a+b-A} - \frac{a}{a+b} = \frac{aA}{(a+b)(a+b-A)} \leq \frac{A}{a+b-A} \leq \frac{A}{B-A} .$$

$\square$

We now lift Lemma 13 to a more general statement.

**Lemma 16** *Let $G = (V, D)$ be a directed graph, with vertex set $V = \{1, \ldots, K\}$, and arc set $D$. Let $\alpha$ be the independence number of $G$, $R \subseteq V$ be a dominating set for $G$ of size $r = |R|$, and $p_1, \ldots, p_K$ be a probability distribution defined over $V$, such that $p_i \geq \beta > 0$, for $i \in R$. Then*

$$\sum_{i=1}^{K} \frac{p_i}{p_i + \sum_{j\,:\,j \to i}\ p_j} \leq 2\alpha \ln\left(1 + \frac{\lceil \frac{K^2}{r\beta} \rceil + K}{\alpha}\right) + 2r .$$

**Proof.** The idea is to appropriately discretize the probability values $p_i$, and then upper bound the discretized counterpart of $\sum_{i=1}^{K} \frac{p_i}{p_i + \sum_{j\,:\,j \to i}\ p_j}$ by reducing to an expression that can be handled by Lemma 13. In order to make this discretization effective, we need to single out the terms $\frac{p_i}{p_i + \sum_{j\,:\,j \to i}\ p_j}$ corresponding to nodes $i \in R$. We first write

$$\begin{aligned}
\sum_{i=1}^{K} \frac{p_i}{p_i + \sum_{j\,:\,j \to i}\ p_j} &= \sum_{i \in R} \frac{p_i}{p_i + \sum_{j\,:\,j \to i}\ p_j} + \sum_{i \notin R} \frac{p_i}{p_i + \sum_{j\,:\,j \to i}\ p_j} \\
&\leq r + \sum_{i \notin R} \frac{p_i}{p_i + \sum_{j\,:\,j \to i}\ p_j}
\end{aligned} \tag{13}$$

and then focus on (13).

Let us discretize the unit interval[12] $(0,1]$ into subintervals $\left(\frac{j-1}{M}, \frac{j}{M}\right]$, $j = 1, \ldots, M$, where $M = \lceil \frac{K^2}{r\beta} \rceil$. Let $\widehat{p}_i = j/M$ be the discretized version of $p_i$, being $j$ the unique integer such that $\widehat{p}_i - 1/M < p_i \leq \widehat{p}_i$. We focus on a single node $i \notin R$ with indegree $d_i^-$. Introduce the shorthand notations $P_i = \sum_{j:j\to i} p_j$ and $\widehat{P}_i = \sum_{j:j\to i} \widehat{p}_j$. We have that $\widehat{P}_i \geq P_i \geq \beta$, since $i$ is dominated by some node $j \in R \cap N_i^-$ such that $p_j \geq \beta$. Moreover, $P_i > \widehat{P}_i - \frac{d_i^-}{M} \geq \beta - \frac{d_i^-}{M} > 0$, and $\widehat{p}_i + \widehat{P}_i \geq \beta$. Hence, for any fixed node $i \notin R$, we can write

$$
\begin{aligned}
\frac{p_i}{p_i + P_i} &\leq \frac{\widehat{p}_i}{\widehat{p}_i + P_i} \\
&< \frac{\widehat{p}_i}{\widehat{p}_i + \widehat{P}_i - \frac{d_i^-}{M}} \\
&\leq \frac{\widehat{p}_i}{\widehat{p}_i + \widehat{P}_i} + \frac{d_i^-/M}{\beta - d_i^-/M} \\
&= \frac{\widehat{p}_i}{\widehat{p}_i + \widehat{P}_i} + \frac{d_i^-}{\beta M - d_i^-} \\
&< \frac{\widehat{p}_i}{\widehat{p}_i + \widehat{P}_i} + \frac{r}{K - r}
\end{aligned}
$$

where in the second-last inequality we used Lemma 15 with $a = \widehat{p}_i$, $b = \widehat{P}_i$, $A = d_i^-/M$, and $B = \beta > d_i^-/M$. Recalling (13), and summing over $i$ then gives

$$
\sum_{i=1}^{K} \frac{p_i}{p_i + P_i} \leq r + \sum_{i \notin R} \frac{\widehat{p}_i}{\widehat{p}_i + \widehat{P}_i} + r = \sum_{i \notin R} \frac{\widehat{p}_i}{\widehat{p}_i + \widehat{P}_i} + 2r \ . \tag{14}
$$

Therefore, we continue by bounding from above the right-hand side of (14). We first observe that

$$
\sum_{i \notin R} \frac{\widehat{p}_i}{\widehat{p}_i + \widehat{P}_i} = \sum_{i \notin R} \frac{\widehat{s}_i}{\widehat{s}_i + \widehat{S}_i} \qquad \text{and} \qquad \widehat{S}_i = \sum_{j:j\to i} \widehat{s}_j \tag{15}
$$

where $\widehat{s}_i = M\widehat{p}_i$, $i = 1, \ldots, K$, are integers. Based on the original graph $G$, we construct a new graph $\widehat{G}$ made up of connected cliques. In particular:

- Each node $i$ of $G$ is replaced in $\widehat{G}$ by a clique $C_i$ of size $\widehat{s}_i$; nodes within $C_i$ are connected by length-two cycles.

- If arc $(i, j)$ is in $G$, then for each node of $C_i$ draw an arc towards each node of $C_j$.

We would like to apply Lemma 13 to $\widehat{G}$. Note that, by the above construction:

- The independence number of $\widehat{G}$ is the same as that of $G$;

- The indegree $\widehat{d_k^-}$ of each node $k$ in clique $C_i$ satisfies $\widehat{d_k^-} = \widehat{s}_i - 1 + \widehat{S}_i$.

---

[12] The zero value is not of our concern here, because if $p_i = 0$, then the corresponding term in (13) can be disregarded.

- The total number of nodes of $\widehat{G}$ is

$$\sum_{i=1}^{K} \widehat{s}_i = M \sum_{i=1}^{K} \widehat{p}_i < M \sum_{i=1}^{K} \left( p_i + \frac{1}{M} \right) = M + K .$$

Hence, we can apply Lemma 13 to $\widehat{G}$ with indegrees $\widehat{d}_k^{-}$, and find that

$$\sum_{i \notin R} \frac{\widehat{s}_i}{\widehat{s}_i + \widehat{S}_i} = \sum_{i \notin R} \sum_{k \in C_i} \frac{1}{1 + \widehat{d}_k^{-}} \leq \sum_{i=1}^{K} \sum_{k \in C_i} \frac{1}{1 + \widehat{d}_k^{-}} \leq 2\alpha \ln \left( 1 + \frac{M+K}{\alpha} \right) .$$

Putting together (14) and (15), and recalling the value of $M$ gives the claimed result. $\qquad\square$

### Proof of Theorem 8

We are now ready to derive the proof of the theorem. We start from the upper bound (4) in the statement of Lemma 7. We want to bound the quantities $|R_t|$ and $Q_t^{(b_t)}$ occurring therein at any step $t$ in which a restart does not occur —the regret for the time steps when a restart occurs is already accounted for by the term $\mathcal{O}\big((\ln K)\ln(KT)\big)$ in (4). Now, Lemma 14 gives

$$|R_t| = \mathcal{O}\big(\alpha(G_t) \ln K\big) .$$

If $\gamma_t = \gamma_t^{(b_t)}$ for any time $t$ when a restart does not occur, it is not hard to see that $\gamma_t = \Omega\big(\sqrt{(\ln K)/(KT)}\big)$. Moreover, Lemma 16 states that

$$Q_t = \mathcal{O}\big(\alpha(G_t) \ln(K^2/\gamma_t) + |R_t|\big) = \mathcal{O}\big(\alpha(G_t) \ln(K/\gamma_t)\big) .$$

Hence,

$$Q_t = \mathcal{O}\big(\alpha(G_t) \ln(KT)\big).$$

Putting together as in (4) gives the desired result.

## C  Technical lemmas and proofs from Section 4.2

Once again, throughout this appendix $\mathbb{E}_t[\,\cdot\,]$ denotes the conditional expectation $\mathbb{E}_t[\,\cdot \mid I_1, I_2, \ldots, I_{t-1}]$. Moreover, as we did in previous appendices, we first condition on the history $I_1, I_2, \ldots, I_{t-1}$, and then take an expectation with respect to that history.

### C.1  Proof of Theorem 9

The following lemmas are of preliminary importance in order to understand the behavior of the ELP.P algorithm. Recall that for a directed graph $G = (V, D)$, with vertex set $V = \{1, \ldots, K\}$, and arc set $D$, we write $\{j : j \to i\}$ to denote the set of nodes $j$ which are in-neighbors of node $i$, where it is understood that node $i$ is an in-neighbor of itself. Similarly, $\{j : i \to j\}$ is the out-neighborhood of node $i$ where, again, node $i$ is an out-neighbor of itself. Let $\Delta_K$ be the $K$-dimensional probability simplex.

**Lemma 17** *Consider a directed graph $G = (V, D)$, with vertex set $V = \{1, \ldots, K\}$, and arc set $D$. Let $\mathtt{mas}(G)$ be the size of a largest acyclic subgraph of $G$. If $s_1, \ldots, s_K$ is a solution to the linear program*

$$\max_{(s_1,\ldots,s_K)\in\Delta_K} \min_{i\in V} \left( \sum_{j\,:\,j\to i} s_j \right) \tag{16}$$

*then we have*

$$\max_{i\in V} \frac{1}{\sum_{j\,:\,j\to i} s_j} \le \mathtt{mas}(G) \ .$$

**Proof.** We first show that the above inequality holds when the right-hand side is replaced by $\gamma(G)$, the domination number of $G$. Let then $R$ be a smallest (minimal) dominating set of $G$, so that $|R| = \gamma(G)$. Consider the valid assignment $s_i = \mathbb{I}\{i \in R\}/\gamma(G)$ for all $i \in V$. This implies that for all $i$, $\sum_{j\,:\,j\to i} s_j \ge 1/\gamma(G)$, because any $i \in V$ is either in $R$ or is dominated by a node in $R$. Therefore, for this particular assignment, we have

$$\max_{i\in V} \frac{1}{\sum_{j\,:\,j\to i} s_j} \le \gamma(G) \ .$$

The assignment returned by the linear program might be different, but it can only make the left-hand side above smaller,[13] so the inequality still holds. Finally, $\gamma(G) \le \mathtt{mas}(G)$ because any set $M \subseteq V$ of nodes belonging to a maximal acyclic subgraph of $G$ is itself a dominating set for $G$. In fact, assuming the contrary, let $j$ be any node such that $j \notin M$. Then, including $j$ in $M$ would create a cycle (because of the maximality of $M$), implying that $j$ is already dominated by some other node in $M$. □

**Lemma 18** *Consider a directed graph $G = (V, D)$, with vertex set $V = \{1, \ldots, K\}$, and arc set $D$. Let $\mathtt{mas}(G)$ be the size of a largest acyclic subgraph of $G$. Let $(p_1, \ldots, p_K) \in \Delta_K$ and $(s_1, \ldots, s_K) \in \Delta_K$ satisfy*

$$\sum_{i=1}^{K} \frac{p_i}{\sum_{j\,:\,j\to i} p_j} \le \mathtt{mas}(G) \qquad and \qquad \max_{i\in V} \frac{1}{\sum_{j\,:\,j\to i} s_j} \le \mathtt{mas}(G)$$

*with $p_i \ge \gamma s_i$, $i \in V$, for some $\gamma > 0$. Then the following relations hold:*

*1.*

$$\sum_{i=1}^{K} \frac{p_i}{\left( \sum_{j\,:\,j\to i} p_j \right)^2} \le \frac{\mathtt{mas}^2(G)}{\gamma} \ ;$$

*2.*

$$\sum_{i=1}^{K} p_i \sum_{j\,:\,i\to j} \frac{p_j}{\sum_{\ell\,:\,\ell\to j} p_\ell} = 1 \ ;$$

---

[13] This can be seen by noting that (16) is equivalent to

$$\min_{(s_1,\ldots,s_K)\in\Delta_K} \max_{i\in V} \frac{1}{\sum_{j\,:\,j\to i} s_j}$$

*3.*
$$\sum_{i=1}^{K} p_i \sum_{j\,:\,i\to j} \frac{p_j}{\left(\sum_{\ell\,:\,\ell\to j} p_\ell\right)^2} \le \texttt{mas}(G)\ ;$$

*4.*
$$\sum_{i=1}^{K} p_i \left(\sum_{j\,:\,i\to j} \frac{p_j}{\sum_{\ell\,:\,\ell\to j} p_\ell}\right)^2 \le \texttt{mas}(G)\ ;$$

*5.*
$$\sum_{i=1}^{K} p_i \left(\sum_{j\,:\,i\to j} \frac{p_j}{\left(\sum_{\ell\,:\,\ell\to j} p_\ell\right)^2}\right)^2 \le \frac{\texttt{mas}^3(G)}{\gamma}\ .$$

**Proof.** Let us introduce the shorthand $q_i = \sum_{j\,:\,j\to i} p_j$, for $i \in V$.

1. We apply Hölder's inequality, and the assumptions of this lemma to obtain

$$
\begin{aligned}
\sum_{i=1}^{K} \frac{p_i}{q_i^2} &= \sum_{i=1}^{K} \left(\frac{p_i}{q_i}\right)\left(\frac{1}{q_i}\right) \\
&\le \left(\sum_{i=1}^{K} \frac{p_i}{q_i}\right)\left(\max_{i\in V} \frac{1}{q_i}\right) \\
&= \left(\sum_{i=1}^{K} \frac{p_i}{q_i}\right)\left(\max_{i\in V} \frac{1}{\sum_{j\,:\,j\to i} p_j}\right) \\
&\le \texttt{mas}(G) \max_{i\in V} \frac{1}{\gamma\left(\sum_{j\,:\,j\to i} s_j\right)} \\
&\le \frac{\texttt{mas}^2(G)}{\gamma}\ .
\end{aligned}
$$

2. We have
$$\sum_{i=1}^{K} \sum_{j\,:\,i\to j} \frac{p_i\,p_j}{q_j} = \sum_{j=1}^{K} \frac{p_j\,q_j}{q_j} = \sum_{j=1}^{K} p_j = 1\ .$$

3. Similar to the previous item, we can write
$$\sum_{i=1}^{K} \sum_{j\,:\,i\to j} \frac{p_i\,p_j}{q_j^2} = \sum_{j=1}^{K} \frac{p_j\,q_j}{q_j^2} = \sum_{j=1}^{K} \frac{p_j}{q_j} \le \texttt{mas}(G)\ .$$

4. From Item 2, and the assumptions of this lemma, we can write

$$
\begin{aligned}
\sum_{i=1}^{K} p_i \left( \sum_{j\,:\,i\to j} \frac{p_j}{q_j} \right)^2 &= \sum_{i=1}^{K} \left( p_i \sum_{j\,:\,i\to j} \frac{p_j}{q_j} \right) \left( \sum_{j\,:\,i\to j} \frac{p_j}{q_j} \right) \\
&\leq \left( \sum_{i=1}^{K} p_i \sum_{j\,:\,i\to j} \frac{p_j}{q_j} \right) \left( \max_{i\in V} \sum_{j\,:\,i\to j} \frac{p_j}{q_j} \right) \\
&\leq \left( \sum_{i=1}^{K} p_i \sum_{j\,:\,i\to j} \frac{p_j}{q_j} \right) \left( \sum_{j=1}^{K} \frac{p_j}{q_j} \right) \\
&= \sum_{j=1}^{K} \frac{p_j}{q_j} \leq \mathtt{mas}(G) \ .
\end{aligned}
$$

5. From Item 1 and Item 3 above, we can write

$$
\begin{aligned}
\sum_{i=1}^{K} p_i \left( \sum_{j\,:\,i\to j} \frac{p_j}{q_j^2} \right)^2 &= \sum_{i=1}^{K} \left( p_i \sum_{j\,:\,i\to j} \frac{p_j}{q_j^2} \right) \left( \sum_{j\,:\,i\to j} \frac{p_j}{q_j^2} \right) \\
&\leq \left( \sum_{i=1}^{K} p_i \sum_{j\,:\,i\to j} \frac{p_j}{q_j^2} \right) \left( \max_{i\in V} \sum_{j\,:\,i\to j} \frac{p_j}{q_j^2} \right) \\
&\leq \left( \sum_{i=1}^{K} p_i \sum_{j\,:\,i\to j} \frac{p_j}{q_j^2} \right) \left( \sum_{i=1}^{K} \frac{p_i}{q_i^2} \right) \\
&\leq \mathtt{mas}(G) \, \frac{\mathtt{mas}^2(G)}{\gamma} \\
&= \frac{\mathtt{mas}^3(G)}{\gamma}
\end{aligned}
$$

concluding the proof. $\qquad\square$

Lemma 18 applies, in particular, to the distributions $s_t = (s_{1,t}, \ldots, s_{K,t})$ and $p_t = (p_{1,t}, \ldots, p_{K,t})$ computed by ELP.P at round $t$. The condition for $p_t$ follows from Lemma 10, while the condition for $s_t$ follows from Lemma 17. In other words, putting together Lemma 10 and Lemma 17 establishes the following corollary.

**Corollary 19** *Let $p_t = (p_{1,t}, \ldots, p_{K,t}) \in \Delta_K$ and $s_t = (s_{1,t}, \ldots, s_{K,t}) \in \Delta_K$ be the distributions generated by ELP.P at round $t$. Then,*

$$
\sum_{i=1}^{K} \frac{p_{i,t}}{\sum_{j\,:\,j\xrightarrow{t} i} p_{j,t}} \leq \mathtt{mas}(G) \qquad \text{and} \qquad \max_{i\in V} \frac{1}{\sum_{j\,:\,j\xrightarrow{t} i} s_{j,t}} \leq \mathtt{mas}(G) \ ,
$$

*with $p_{i,t} \geq \gamma_t \, s_{i,t}$, for all $i = 1, \ldots, K$.*

34

For the next result, we need the following version of Freedman's inequality [14] (see also [9, Lemma A.8]).

**Lemma 20** *Let $X_1, \ldots, X_T$ be a martingale difference sequence with respect to the filtration $\{\mathcal{F}_t\}_{t=1,\ldots,T}$, and with $|X_i| \leq B$ almost surely for all $i$. Also, let $V > 0$ be a fixed upper bound on $\sum_{t=1}^T \mathbb{E}\big[X_t^2 \mid \mathcal{F}_{t-1}\big]$. Then for any $\delta \in (0,1)$, it holds with probability at least $1 - \delta$*

$$\sum_{t=1}^T X_t \leq \sqrt{2 \ln\left(\frac{1}{\delta}\right) V} + \frac{B}{2} \ln\left(\frac{1}{\delta}\right) .$$

**Lemma 21** *Let $\{a_t\}_{t=1}^T$ be an arbitrary sequence of positive numbers, and let $s_t = (s_{1,t}, \ldots, s_{K,t})$ and $p_t = (p_{1,t}, \ldots, p_{K,t})$ be the probability distributions computed by ELP.P at the $t$-th round. Then, with probability at least $1 - \delta$,*

$$\sum_{t=1}^T \sum_{i=1}^K a_t p_{i,t}(\widehat{g}_{i,t} - g_{i,t})$$

$$\leq \sqrt{2 \ln\left(\frac{1}{\delta}\right) \sum_{t=1}^T a_t^2 \, \mathtt{mas}(G_t)} + \frac{1}{2} \ln\left(\frac{1}{\delta}\right) \max_{t=1\ldots T} \big(a_t \, \mathtt{mas}(G_t)\big) + \beta \sum_{t=1}^T a_t \, \mathtt{mas}(G_t) . \quad (17)$$

**Proof.** Recall that $q_{i,t} = \sum_{j\,:\,j \xrightarrow{t} i} p_{j,t}$, for $i \in V$, and let

$$\tilde{g}_{i,t} = \frac{g_{i,t} \, \mathbb{I}\{i \in S_{I_t,t}\}}{q_{i,t}}$$

with $g_{i,t} = 1 - \ell_{i,t}$. Note that $\widehat{g}_{i,t}$ in Figure 3 satisfies $\widehat{g}_{i,t} = \tilde{g}_{i,t} + \frac{\beta}{q_{i,t}}$, so that we can upper bound the left-hand side of (17) by

$$\sum_{t=1}^T \sum_{i=1}^K a_t p_{i,t}(\tilde{g}_{i,t} - g_{i,t}) + \beta \sum_{t=1}^T \sum_{i=1}^K \frac{a_t p_{i,t}}{q_{i,t}}$$

which by Corollary 19 is at most

$$\sum_{t=1}^T \sum_{i=1}^K a_t p_{j,t}(\tilde{g}_{i,t} - g_{i,t}) + \beta \sum_{t=1}^T a_t \, \mathtt{mas}(G_t) . \quad (18)$$

It is easy to verify that $\sum_{i=1}^K a_t p_{i,t}(\tilde{g}_{i,t} - g_{i,t})$ is a martingale difference sequence (indexed by $t$), since $\tilde{g}_{i,t}$ is an unbiased estimate of $g_{i,t}$ conditioned on the previous rounds. Moreover,

$$\sum_{i=1}^K a_t p_{i,t}(\tilde{g}_{i,t} - g_{i,t}) = \sum_{i=1}^K a_t p_{i,t}\Big(\frac{\mathbb{I}\{i \in S_{I_t,t}\}}{q_{i,t}} - 1\Big) g_{i,t} \leq a_t \sum_{i=1}^K \frac{p_{i,t}}{q_{i,t}} \leq \max_{t=1,\ldots,T} a_t \, \mathtt{mas}(G_t)$$

35

and

$$\mathbb{E}_t\left[\left(\sum_{i=1}^{K}a_t p_{i,t}(\tilde{g}_{i,t}-g_{i,t})\right)^2\right] \leq a_t^2\,\mathbb{E}_t\left[\left(\sum_{i=1}^{K}p_{i,t}\tilde{g}_{i,t}\right)^2\right]$$

$$\leq a_t^2\sum_{i=1}^{K}p_{i,t}\left(\sum_{j\,:\,i\xrightarrow{t}j}p_{j,t}\frac{1}{q_{j,t}}\right)^2$$

$$\leq a_t^2\,\mathtt{mas}(G_t)$$

by Lemma 18, item 4. Therefore, by invoking Lemma 20, we get that with probability at least $1-\delta$,

$$\sum_{t=1}^{T}\sum_{j=1}^{K}a_t p_{j,t}(\tilde{g}_{j,t}-g_{j,t}) \leq \sqrt{2\ln\left(\frac{1}{\delta}\right)\sum_{t=1}^{T}a_t^2\,\mathtt{mas}(G_t)} + \frac{1}{2}\ln\left(\frac{1}{\delta}\right)\max_{t=1,\dots,T}a_t\,\mathtt{mas}(G_t)\,.$$

Substituting into Eq. (18), the lemma follows. $\qquad\square$

**Lemma 22** *Let $s_t = (s_{1,t},\dots,s_{K,t})$ and $p_t = (p_{1,t},\dots,p_{K,t})$ be the probability distributions computed by ELP.P, run with $\beta \leq 1/4$, at round $t$. Then, with probability at least $1-\delta$,*

$$\sum_{t=1}^{T}\sum_{i=1}^{K}p_{i,t}\widehat{g}_{i,t}^2 \leq \sum_{t=1}^{T}\left(\frac{\beta^2\,\mathtt{mas}^2(G_t)}{\gamma_t}+2\,\mathtt{mas}(G_t)\right) + \sqrt{2\ln\left(\frac{1}{\delta}\right)\sum_{t=1}^{T}\left(\frac{4\beta^2\,\mathtt{mas}^4(G_t)}{\gamma_t^2}+\frac{3\,\mathtt{mas}^3(G_t)}{\gamma_t}\right)}$$

$$+\ln\left(\frac{1}{\delta}\right)\max_{t=1,\dots,T}\frac{\mathtt{mas}^2(G_t)}{\gamma_t}\,.$$

**Proof.** Recall that $q_{i,t}=\sum_{j\,:\,j\xrightarrow{t}i}p_{j,t}$, for $i\in V$. By the way we defined $\widehat{g}_{i,t}$ and Lemma 18, item 1, we have that

$$\sum_{i=1}^{K}p_{i,t}\widehat{g}_{i,t}^2 \leq \sum_{i=1}^{K}p_{i,t}\left(\frac{1+\beta}{q_{i,t}}\right)^2 \leq \frac{(1+\beta)^2\mathtt{mas}^2(G_t)}{\gamma_t}\,.$$

Moreover, from $g_{i,t}\leq 1$, and again using Lemma 18, item 1, we can write

$$\mathbb{E}_t\left[\left(\sum_{j=1}^{K}p_{j,t}\widehat{g}_{j,t}^2\right)^2\right] \leq \sum_{i=1}^{K}p_{i,t}\left(\sum_{j=1}^{K}\frac{p_{j,t}}{(q_{j,t})^2}\left(\mathbb{I}\{i\xrightarrow{t}j\}+\beta\right)^2\right)^2$$

$$=\sum_{i=1}^{K}p_{i,t}\left(\beta^2\sum_{j=1}^{K}\frac{p_{j,t}}{(q_{j,t})^2}+(2\beta+1)\sum_{j\,:\,i\xrightarrow{t}j}\frac{p_{j,t}}{(q_{j,t})^2}\right)^2$$

$$\leq\sum_{i=1}^{K}p_{i,t}\left(\frac{\beta^2\,\mathtt{mas}^2(G_t)}{\gamma_t}+(2\beta+1)\sum_{j\,:\,i\xrightarrow{t}j}\frac{p_{j,t}}{(q_{j,t})^2}\right)^2$$

36

which by expanding, using Lemma 18, items 3 and 5, and slightly simplifying, is at most

$$\frac{(\beta^4 + 2\beta^2(2\beta + 1))\,\mathtt{mas}^4(G_t)}{\gamma_t^2} + \frac{(2\beta + 1)^2\,\mathtt{mas}^3(G_t)}{\gamma_t} \leq \frac{4\beta^2\,\mathtt{mas}^4(G_t)}{\gamma_t^2} + \frac{3\,\mathtt{mas}^3(G_t)}{\gamma_t}$$

the last inequality exploiting the assumption $\beta \leq 1/4$. Invoking Lemma 20 we get that with probability at least $1 - \delta$

$$\sum_{t=1}^{T}\sum_{i=1}^{K} p_{i,t}\widehat{g}_{i,t}^2 \leq \sum_{t=1}^{T}\sum_{i=1}^{K} p_{i,t}\mathbb{E}_t[\widehat{g}_{i,t}^2] + \sqrt{2\ln\left(\frac{1}{\delta}\right)\sum_{t=1}^{T}\left(\frac{4\beta^2\,\mathtt{mas}^4(G_t)}{\gamma_t^2} + \frac{3\,\mathtt{mas}^3(G_t)}{\gamma_t}\right)}$$

$$+ \frac{(1+\beta)^2}{2}\ln\left(\frac{1}{\delta}\right)\max_{t=1...T}\frac{\mathtt{mas}^2(G_t)}{\gamma_t} \quad . \tag{19}$$

Finally, from $g_{i,t} \leq 1$, Lemma 18, item 1, and the assumptions of this lemma, we have

$$\sum_{i=1}^{K} p_{i,t}\mathbb{E}_t[\widehat{g}_{i,t}^2] \leq \sum_{i=1}^{K} p_{i,t}\sum_{j=1}^{K} p_{j,t}\left(\frac{\mathbb{I}\{j \xrightarrow{t} i\} + \beta}{q_{i,t}}\right)^2$$

$$= \beta^2\sum_{i=1}^{K}\frac{p_{i,t}}{(q_{i,t})^2} + (2\beta + 1)\sum_{i=1}^{K} p_{i,t}\sum_{j\,:\,j\xrightarrow{t}i}\frac{p_{j,t}}{(q_{i,t})^2}$$

$$= \beta^2\sum_{i=1}^{K}\frac{p_{i,t}}{(q_{i,t})^2} + (2\beta + 1)\sum_{i=1}^{K}\frac{p_{i,t}}{q_{i,t}}$$

$$\leq \frac{\beta^2\,\mathtt{mas}^2(G_t)}{\gamma_t} + (2\beta + 1)\,\mathtt{mas}(G_t)$$

$$\leq \frac{\beta^2\,\mathtt{mas}^2(G_t)}{\gamma_t} + 2\,\mathtt{mas}(G_t)$$

where we used again $\beta \leq 1/4$. Plugging this back into Eq. (19) the result follows. $\qquad\square$

**Lemma 23** *Suppose that the ELP.P algorithm is run with $\beta \leq 1/4$. Then it holds with probability at least that $1 - \delta$ that for any $i = 1, \dots, K$,*

$$\sum_{t=1}^{T}\widehat{g}_{i,t} \geq \sum_{t=1}^{T} g_{i,t} - \frac{\ln(K/\delta)}{\beta} \quad .$$

**Proof.** Let $\lambda > 0$ be a parameter to be specified later, and recall that $\mathbb{E}_t$ denotes the expectation at round $t$, conditioned on all previous rounds. Since $\exp(x) \leq 1 + x + x^2$ for $x \leq 1$, we have by

definition of $\widehat{g}_{i,t}$ that

$$\mathbb{E}_t\left[\exp\left(\lambda(g_{i,t}-\widehat{g}_{i,t})\right)\right] = \mathbb{E}_t\left[\exp\left(\lambda\left(g_{i,t}-\frac{g_{i,t}\mathbb{I}\{I_t\xrightarrow{t}i\}}{q_{i,t}}\right)-\frac{\beta\lambda}{q_{i,t}}\right)\right]$$

$$\leq \left(1+\mathbb{E}_t\left[\lambda\left(g_{i,t}-\frac{g_{i,t}\mathbb{I}\{I_t\xrightarrow{t}i\}}{q_{i,t}}\right)\right]+\mathbb{E}_t\left[\left(\lambda\left(g_{i,t}-\frac{g_{i,t}\mathbb{I}\{I_t\xrightarrow{t}i\}}{q_{i,t}}\right)\right)^2\right]\right)\exp\left(-\frac{\beta\lambda}{q_{i,t}}\right)$$

$$\leq \left(1+0+\lambda^2\mathbb{E}_t\left[\left(\frac{g_{i,t}\mathbb{I}\{I_t\xrightarrow{t}i\}}{q_{i,t}}\right)^2\right]\right)\exp\left(-\frac{\beta\lambda}{q_{i,t}}\right)$$

$$\leq \left(1+\lambda^2\sum_{j:j\xrightarrow{t}i}\frac{p_{j,t}}{(q_{i,t})^2}\right)\exp\left(-\frac{\beta\lambda}{q_{i,t}}\right)$$

$$= \left(1+\frac{\lambda^2}{q_{i,t}}\right)\exp\left(-\frac{\beta\lambda}{q_{i,t}}\right)\ .$$

Picking $\lambda = \beta$, and using the fact that $(1+x)\exp(-x) \leq 1$, we get that this expression is at most 1. As a result, we have

$$\mathbb{E}\left[\exp\left(\lambda\sum_{t=1}^{T}(g_{i,t}-\widehat{g}_{i,t})\right)\right] \leq 1\ .$$

Now, by a standard Chernoff technique, we know that for any $\lambda > 0$,

$$\mathbb{P}\left(\sum_{t=1}^{T}(g_{i,t}-\widehat{g}_{i,t}) > \epsilon\right) \leq \exp(-\lambda\epsilon)\mathbb{E}\left[\exp\left(\lambda\sum_{t=1}^{T}(g_{i,t}-\widehat{g}_{i,t})\right)\right]\ .$$

In particular, for our choice of $\lambda$, we get the bound

$$\mathbb{P}\left(\sum_{t=1}^{T}(g_{i,t}-\widehat{g}_{i,t}) > \epsilon\right) \leq \exp\left(-\beta\epsilon\right)\ .$$

Substituting $\delta = \exp(-\beta\epsilon)$, solving for $\epsilon$, and using a union bound to make the result hold simultaneously for all $i$, the result follows. $\square$

**Proof of Theorem 9**

With these key lemmas at hand, we can now turn to prove Theorem 9. We have

$$\frac{W_{t+1}}{W_t} = \sum_{i\in V}\frac{w_{i,t+1}}{W_t} = \sum_{i\in V}\frac{w_{i,t}}{W_t}\exp(\eta\widehat{g}_{i,t})\ . \tag{20}$$

Now, by definition of $q_{i,t}$ and $\gamma_t$ in Algorithm 3 we have

$$q_{i,t} \geq \gamma_t\sum_{j:j\xrightarrow{t}i}s_{j,t} \geq (1+\beta)\eta$$

38

for all $i \in V$, so that

$$\eta \widehat{g}_{j,t} \;\le\; \eta \max_{i \in V} \left( \frac{1+\beta}{q_{i,t}} \right) \;\le\; 1 \;.$$

Using the definition of $p_{i,t}$ and the inequality $\exp(x) \le 1 + x + x^2$ for any $x \le 1$, we can then upper bound the right-hand side of (20) by

$$\sum_{i \in V} \frac{p_{i,t} - \gamma_t s_{i,t}}{1 - \gamma_t} \left( 1 + \eta \widehat{g}_{i,t} + \eta^2 \widehat{g}_{i,t}^2 \right) \;\le\; 1 + \frac{\eta}{1 - \gamma_t} \sum_{i \in V} p_{i,t} \widehat{g}_{i,t} + \frac{\eta^2}{1 - \gamma_t} \sum_{i=1}^{K} p_{i,t} \widehat{g}_{i,t}^2 \;.$$

Taking logarithms and using the fact that $\ln(1 + x) \le x$, we get

$$\ln \left( \frac{W_{t+1}}{W_t} \right) \;\le\; \frac{\eta}{1 - \gamma_t} \sum_{i \in V} p_{i,t} \widehat{g}_{i,t} + \frac{\eta^2}{1 - \gamma_t} \sum_{i \in V} p_{i,t} \widehat{g}_{i,t}^2 \;.$$

Summing over all $t$, and canceling the resulting telescopic series, we get

$$\ln \left( \frac{W_{T+1}}{W_1} \right) \;\le\; \sum_{t=1}^{T} \sum_{i \in V} \frac{\eta}{1 - \gamma_t} p_{i,t} \widehat{g}_{i,t} + \sum_{t=1}^{T} \sum_{i \in V} \frac{\eta^2}{1 - \gamma_t} p_{i,t} \widehat{g}_{i,t}^2. \tag{21}$$

Also, for any fixed action $k$, we have

$$\ln \left( \frac{W_{T+1}}{W_1} \right) \ge \ln \left( \frac{w_{j,T+1}}{W_1} \right) = \eta \sum_{t=1}^{T} \widehat{g}_{k,t} - \ln K \;. \tag{22}$$

Combining Eq. (21) with Eq. (22), and slightly rearranging and simplifying, we get

$$\sum_{t=1}^{T} \widehat{g}_{k,t} - \sum_{t=1}^{T} \sum_{i \in V} p_{i,t} \widehat{g}_{i,t}$$

$$\le \frac{\ln K}{\eta} + \frac{\eta}{1 - \max\limits_{t=1,\dots,T} \gamma_t} \sum_{t=1}^{T} \sum_{i \in V} p_{i,t} \widehat{g}_{i,t}^2 + \frac{1}{1 - \max\limits_{t=1,\dots,T} \gamma_t} \sum_{t=1}^{T} \sum_{i \in V} \gamma_t p_{i,t} \widehat{g}_{i,t} \;. \tag{23}$$

We now start to apply the various lemmas, using a union bound. To keep things manageable, we will use asymptotic notation to deal with second-order terms. In particular, we will use $\widetilde{\mathcal{O}}$ notation to hide numerical constants and logarithmic factors[14]. By definition of $\beta$ and $\gamma_t$, as well as Corollary 19, it is easy to verify[15] that

$$\beta = \widetilde{\mathcal{O}}(\eta) \qquad \gamma_t = \widetilde{\mathcal{O}}(\eta \, \mathtt{mas}(G_t)) \qquad \gamma_t \in \left[ \eta, \frac{1}{2} \right] \;. \tag{24}$$

---

[14] Technically, $\widetilde{\mathcal{O}}(f) = O(f \log^{O(1)} f)$. In our $\widetilde{\mathcal{O}}$ we also ignore factors that depend logarithmically on $K$ and $1/\delta$.
[15] The bound for $\beta$ is by definition. The bound for $\gamma_t$ holds since by Lemma 17 and the assumptions that $\eta \le 1/(3K)$ and $\beta \le 1/4$ we have

$$\gamma_t = \frac{(1+\beta)\eta}{\min_{i \in V} \sum_{j \,:\, j \xrightarrow{t} i} s_{j,t}} \le (1+\beta)\eta \, \mathtt{mas}(G_t) \le \frac{(1+\beta) \max(G_t)}{3K} \le \frac{1 + 1/4}{3} < 1/2 \;.$$

First, by Lemma 21, we have with probability at least $1 - \delta$ that

$$\sum_{t=1}^{T}\sum_{i=1}^{K}p_{i,t}\widehat{g}_{i,t} \leq \sum_{t=1}^{T}\sum_{i=1}^{K}p_{i,t}g_{i,t} + \sqrt{2\ln\left(\frac{1}{\delta}\right)\sum_{t=1}^{T}\mathtt{mas}(G_t)} + \beta\sum_{t=1}^{T}\mathtt{mas}(G_t) + \widetilde{\mathcal{O}}\left(\max_{t=1\ldots T}\mathtt{mas}(G_t)\right) .$$

(25)

Moreover, by Azuma's inequality, we have with probability at least $1 - \delta$ that

$$\sum_{t=1}^{T}\sum_{i=1}^{K}p_{i,t}g_{i,t} \leq \sum_{t=1}^{T}g_{I_t,t} + \sqrt{\frac{\ln(1/\delta)}{2}T} .$$

(26)

Second, again by Lemma 21 and the conditions (24), we have with probability at least $1 - \delta$ that

$$
\begin{aligned}
\sum_{t=1}^{T}\sum_{i=1}^{K}\gamma_t p_{i,t}\widehat{g}_{i,t} &\leq \sum_{t=1}^{T}\sum_{i=1}^{K}\gamma_t p_{i,t}g_{i,t} + \widetilde{\mathcal{O}}\left(\max_{t=1,\ldots,T}\mathtt{mas}^2(G_t)(1 + \sqrt{T}\eta + T\eta^2)\right) \\
&\leq \sum_{t=1}^{T}\gamma_t + \widetilde{\mathcal{O}}\left(\max_{t=1,\ldots,T}\mathtt{mas}^2(G_t)(1 + \sqrt{T}\eta + T\eta^2)\right) .
\end{aligned}
$$

(27)

Third, by Lemma 22, and conditions (24), we have with probability at least $1 - \delta$ that for all $i$,

$$\sum_{t=1}^{T}p_{i,t}\widehat{g}_{i,t}^2 \leq 2\sum_{t=1}^{T}\mathtt{mas}(G_t) + \left(\max_{t=1,\ldots,T}(\mathtt{mas}^2(G_t))\right)\widetilde{\mathcal{O}}\left(T\eta + \frac{1}{\eta} + \sqrt{T\left(1 + \frac{1}{\eta}\right)}\right) .$$

(28)

Fourth, by Lemma 23, we have with probability at least $1 - \delta$ that

$$\sum_{t=1}^{T}\widehat{g}_{k,t} \geq \sum_{t=1}^{T}g_{k,t} - \frac{\ln(K/\delta)}{\beta} .$$

(29)

Combining Eq. (25), Eq. (26), Eq. (27), Eq. (28) and Eq. (29) with a union bound (i.e., replacing $\delta$ by $\delta/5$), substituting back into Eq. (23), and slightly simplifying, we get that with probability at least $1 - \delta$, $\sum_{t=1}^{T}(g_{k,t} - g_{I_t,t})$ is at most

$$\sqrt{2\ln\left(\frac{5}{\delta}\right)\sum_{t=1}^{T}\mathtt{mas}(G_t)} + \beta\sum_{t=1}^{T}\mathtt{mas}(G_t) + \sqrt{\frac{\ln(5/\delta)}{2}T} + \frac{\ln(5K/\delta)}{\beta} + \frac{\ln K}{\eta}$$

$$+ 4\eta\sum_{t=1}^{T}\mathtt{mas}(G_t) + 2\sum_{t=1}^{T}\gamma_t + (1 + \sqrt{T}\eta + T\eta^2)\widetilde{\mathcal{O}}\left(\max_{t=1,\ldots,T}(\mathtt{mas}^2(G_t))\right) .$$

Substituting in the values of $\beta$ and $\gamma_t$, overapproximating, and simplifying (in particular, using Corollary 19 to upper bound $\gamma_t$ by $(1 + \beta)\eta\,\mathtt{mas}(G_t)$), we get the upper bound

$$\sqrt{5\ln\left(\frac{5}{\delta}\right)\sum_{t=1}^{T}\mathtt{mas}(G_t)} + \frac{2\ln(5K/\delta)}{\eta} + 12\eta\sqrt{\frac{\ln(5K/\delta)}{\ln K}}\sum_{t=1}^{T}\mathtt{mas}(G_t)$$

$$+ \widetilde{\mathcal{O}}(1 + \sqrt{T}\eta + T\eta^2)\left(\max_{t=1,\ldots,T}(\mathtt{mas}^2(G_t))\right) .$$

In particular, by picking $\eta$ such that

$$\eta^2 = \frac{1}{6} \frac{\sqrt{\ln(5K/\delta)\,(\ln K)}}{\sum_{t=1}^{T} m_t}$$

noting that this implies $\eta = \widetilde{\mathcal{O}}(1/\sqrt{T})$, and overapproximating once more, we get the second bound

$$\sum_{t=1}^{T} (g_{k,t} - g_{I_t,t}) \leq 10 \frac{\ln^{1/4}(5K/\delta)}{\ln^{1/4} K} \sqrt{\ln\left(\frac{5K}{\delta}\right) \sum_{t=1}^{T} m_t} + \widetilde{\mathcal{O}}(T^{1/4}) \left( \max_{t=1,\dots,T} \mathtt{mas}^2(G_t) \right) .$$

To conclude, we simply plug in $\ell_{i,t} = 1 - g_{i,t}$ for all $i$ and $t$, thereby obtaining the claimed results.