# Path Planning Problems with Side Observations—When Colonels Play Hide-and-Seek

Dong Quan Vu, Patrick Loiseau, Alonso Silva, Long Tran-Thanh

HAL Id: hal-02375789

https://inria.hal.science/hal-02375789v2

Submitted on 15 Mar 2021

# Path Planning Problems with Side Observations—When Colonels Play Hide-and-Seek

**Dong Quan Vu,**[1] **Patrick Loiseau,**[2] **Alonso Silva,**[3] **Long Tran-Thanh**[4]

[1]Nokia Bell Labs France, AAAID Department, [2]Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LIG & MPI-SWS,
[3]Safran Tech, Signal and Information Technologies, [4]Univ. of Southampton, School of Electronics and Computer Science
quan_dong.vu@nokia.com, patrick.loiseau@inria.fr, alonso.silva-allende@safrangroup.com, l.tran-thanh@soton.ac.uk

## Abstract

Resource allocation games such as the famous Colonel Blotto (CB) and Hide-and-Seek (HS) games are often used to model a large variety of practical problems, but only in their one-shot versions. Indeed, due to their extremely large strategy space, it remains an open question how one can efficiently learn in these games. In this work, we show that the online CB and HS games can be cast as path planning problems with side-observations (SOPPP): at each stage, a learner chooses a path on a directed acyclic graph and suffers the sum of losses that are adversarially assigned to the corresponding edges; and she then receives semi-bandit feedback with side-observations (i.e., she observes the losses on the chosen edges plus some others). We propose a novel algorithm, EXP3-OE, the first-of-its-kind with guaranteed efficient running time for SOPPP without requiring any auxiliary oracle. We provide an expected-regret bound of EXP3-OE in SOPPP matching the order of the best benchmark in the literature. Moreover, we introduce additional assumptions on the observability model under which we can further improve the regret bounds of EXP3-OE. We illustrate the benefit of using EXP3-OE in SOPPP by applying it to the online CB and HS games.

## 1 Introduction

Resource allocation games have been studied profoundly in the literature and showed to be very useful to model many practical situations, including online decision problems, see e.g. (Blocki et al. 2013; Bower and Gilbert 2005; Korzhyk, Conitzer, and Parr 2010; Zhang, Lesser, and Shenoy 2009). In particular, two of the most renowned are the Colonel Blotto game (henceforth, CB game) and the Hide-and-Seek game (henceforth, HS game). In the (one-shot) *CB game*, two players, each with a fixed amount of budget, simultaneously allocate their indivisible resources (often referred to as troops) on $n \in \mathbb{N}$ battlefields, each player's payoff is the aggregate of the values of battlefields where she has a higher allocation. The scope of applications of the CB games includes a variety of problems; for instance, in security where resources correspond to security forces (e.g., (Chia 2012; Schwartz, Loiseau, and Sastry 2014)), in politics where budget are distributed to attract voters (e.g., (Kovenock and Roberson 2012; Roberson 2006)), and in advertising for distributing the ads' broadcasting time (e.g., (Masucci and

Silva 2014; 2015)). On the other hand, in the (one-shot) *HS game*, a seeker chooses $n$ among $k$ locations ($n \leq k$) to search for a hider, who chooses the probability of hiding in each location. The seeker's payoff is the summation of the probability that the hider hides in the chosen locations and the hider's payoff is the probability that she successfully escapes the seeker's pursuit. Several variants of the HS games are used to model surveillance situations (Bhattacharya, Başar, and Falcone 2014), anti-jamming problems (Navda et al. 2007; Wang and Liu 2016), vehicles control (Vidal et al. 2002), etc.

Both the CB and the HS games have a long-standing history (originated by (Borel 1921) and (Von Neumann 1953), respectively); however, the results achieved so-far in these games are mostly limited to their one-shot and full-information version (see e.g., (Behnezhad et al. 2017; Gross and Wagner 1950; Roberson 2006; Schwartz, Loiseau, and Sastry 2014; Vu, Loiseau, and Silva 2018) for CB games and (Hespanha, Prandini, and Sastry 2000; Yavin 1987) for HS games). On the contrary, in most of the applications (e.g., telecommunications, web security, advertising), a more natural setting is to consider the case where the game is played repeatedly and players have access only to incomplete information at each stage. In this setting, players are often required to sequentially learn the game on-the-fly and adjust the trade-off between exploiting known information and exploring to gain new information. Thus, this work focuses on the following sequential learning problems:

($i$) The *online CB game*: fix $k, n \in \mathbb{N}$ ($k, n \geq 1$); at each stage, a learner who has the budget $k$ plays a CB game against some adversaries across $n$ battlefields; at the end of the stage, she receives limited feedback that is the gain (loss) she obtains from each battlefield (but not the adversaries' strategies). The battlefields' values can change over time and they are unknown to the learner before making the decision at each stage. This setting is generic and covers many applications of the CB game. For instance, in radio resource allocation problem (in a cognitive radio network), a solution that balances between efficiency and fairness is to provide the users fictional budgets (the same budget at each stage) and let them bid across $n$ spectrum carriers simultaneously to compete for obtaining as many bandwidth

portions as possible, the highest bidder to each carrier wins the corresponding bandwidth (see e.g., (Chien et al. 2019)). At the end of each stage, each user observes her own data rate (the gain/loss) achieved via each carrier (corresponding to battlefields' values) but does not know other users' bids. Note that the actual data rate can be noisy and change over time. Moreover, users can enter and leave the system so no stochastic assumption shall be made for the adversaries' decisions.

($ii$) The *online HS game*: fix $k, n \in \mathbb{N}$ ($k, n \geq 1$ and $n \leq k$); at each stage, the learner is a seeker who plays the same HS game (with $k$ and $n$) against an adversary; at the end of the stage, the seeker only observes the gains/losses she suffers from the locations she chose. This setting is practical and one of the motivational examples is the spectrum sensing problem in opportunistic spectrum access context (see e.g., (Yucek and Arslan 2009)). At each stage, a secondary user (the learner) chooses to send the sensing signal to at most $n$ among $k$ channels (due to energy constraints, she cannot sense all channels) with the objective of sensing the channels with the availability as high as possible. The leaner can only measure the reliably (the gain/loss) of the channels that she sensed. Note that the channels' availability depend on primary users' decisions that is non-stochastic.

A formal definition of these problems is given in Section 4; hereinafter, we reuse the term CB game and HS game to refer to this sequential learning version of the games. The main challenge here is that the strategy space is exponential in the natural parameters (e.g., number of troops and battlefields in the CB game, number of locations in the HS game); hence how to efficiently learn in these games is an open question.

Our **first contribution** towards solving this open question is to show that the CB and HS games can be cast as a *Path Planning Problem* (henceforth, PPP), one of the most well-studied instances of the *Online Combinatorial Optimization* framework (henceforth, OCOMB; see (Chen, Wang, and Yuan 2013) for a survey). In PPPs, given a directed graph a source and a destination, at each stage, a learner chooses a path from the source to the destination; simultaneously, a loss is adversarially chosen for each edge; then, the learner suffers the aggregate of edges' losses belonging to her chosen path. The learner's goal is to minimize regret. The information that the learner receives in the CB and HS games as described above straightforwardly corresponds to the so-called *semi-bandit* feedback setting of PPPs, i.e., at the end of each stage, the learner observes the edges' losses belonging to her chosen path (see Section 4 for more details). However, the specific structure of the considered games also allows the learner to deduce (without any extra cost) from the semi-bandit feedback the losses of some of the other edges that may not belong to the chosen path; these are called *side-observations*. Henceforth, we will use the term SOPPP to refer to this PPP under semi-bandit feedback with side-observations.

SOPPP is a special case of OCOMB with side-observations (henceforth, SOCOMB) studied by (Kocák et al. 2014) and, following their approach, we will use *obser-*

*vation graphs*[1] (defined in Section 2) to capture the learner's observability. (Kocák et al. 2014) focuses on the class of Follow-the-Perturbed-Leader (FPL) algorithms (originated from (Kalai and Vempala 2005)) and proposes an algorithm named FPL-IX for SOCOMB, which could be applied directly to SOPPP. However, this faces two main problems: ($i$) the efficiency of FPL-IX is only guaranteed with high-probability (as it depends on the geometric sampling technique) and it is still super-linear in terms of the time horizon, thus there is still room for improvements; ($ii$) FPL-IX requires that there exists an efficient oracle that solves an optimization problem at each stage. Both of these issues are incompatible with our goal of learning in the CB and HS games: although the probability that FPL-IX fails to terminate is small, this could lead to issues in implementing it in practice where the learner is obliged to quickly give a decision in each stage; it is unclear which oracle should be used in applying FPL-IX to the CB and HS games.

In this paper, we focus instead on another prominent class of OCOMB algorithms, called EXP3 (Auer et al. 2002; Freund and Schapire 1997). One of the key open questions in this field is how to design a variant of EXP3 with efficient running time and good regret guarantees for OCOMB problems in each feedback setting (see, e.g., (Cesa-Bianchi and Lugosi 2012)). Then, our **second contribution** is to propose an EXP3-type algorithm for SOPPPs that solves both of the aforementioned issues of FPL-IX and provides good regret guarantees; i.e., we give an affirmative answer to an important subset of the above-mentioned open problem. In more details, this contribution is three-fold: ($i$) We propose a *novel algorithm,* EXP3-OE*,* that is applicable to any instance of SOPPP. Importantly, EXP3-OE is always guaranteed to run efficiently (i.e., in polynomial time in terms of the number of edges of the graph in SOPPP) without the need of any auxiliary oracle; ($ii$) We prove that EXP3-OE guarantees an upper-bound on the expected regret matching in order with the best benchmark in the literature (the FPL-IX algorithm). We also prove further improvements under additional assumptions on the observation graphs that have been so-far ignored in the literature; ($iii$) We demonstrate the benefit of using the EXP3-OE algorithm in the CB and HS games.

Note importantly that the SOPPP model (and the EXP3-OE algorithm) can be applied into many problems beyond the CB and HS games, e.g., auctions, recommendation systems. To highlight this and for the sake of conciseness, we first study the generic model of SOPPP in Section 2 and present our second contribution in Section 3, i.e., the EXP3-OE algorithm in SOPPPs; we delay the formal definition of the CB and HS games, together with the analysis on running EXP3-OE in these games (i.e., our first contribution) to Section 4.

---

[1]The observation graphs, proposed by (Kocák et al. 2014) and used here for SOPPP, extend the side-observations model for multi-armed bandits problems studied by (Alon et al. 2015; 2013; Mannor and Shamir 2011). Indeed, they capture side-observations between edges whereas the side-observations model considered by (Alon et al. 2015; 2013; Mannor and Shamir 2011) is between actions, i.e., paths in PPPs.

Throughout the paper, we use bold symbols to denote vectors, e.g., $\boldsymbol{z} \in \mathbb{R}^n$, and $\boldsymbol{z}(i)$ to denote the $i$-th element. For any $m \geq 1$, the set $\{1, 2, \ldots, m\}$ is denoted by $[m]$ and the indicator function of a set $A$ is denoted by $\mathbb{I}_A$. For graphs, we write either $e \in \boldsymbol{p}$ or $\boldsymbol{p} \ni e$ to refer that an edge $e$ belongs to a path $\boldsymbol{p}$. Finally, we use $\tilde{\mathcal{O}}$ as a version of the big-O asymptotic notation that ignores the logarithmic terms.

## 2 Path Planning Problems with Side-Observations (SOPPP) Formulation

As discussed in Section 1, motivated by the CB and HS games, we propose the path planning problem with semi-bandit and side-observations feedback (SOPPP).

**SOPPP model.** Consider a directed acyclic graph (henceforth, DAG), denoted by $G$, whose set of vertices and set of edges are respectively denoted by $\mathcal{V}$ and $\mathcal{E}$. Let $V := |\mathcal{V}| \geq 2$ and $E := |\mathcal{E}| \geq 1$; there are two special vertices, a source and a destination, that are respectively called $s$ and $d$. We denote by $\mathcal{P}$ the set of all *paths* starting from $s$ and ending at $d$; let us define $P := |\mathcal{P}|$. Each path $\boldsymbol{p} \in \mathcal{P}$ corresponds to a vector in $\{0,1\}^E$ (thus, $\mathcal{P} \subset \{0,1\}^E$) where $\boldsymbol{p}(e) = 1$ if and only if edge $e \in \mathcal{E}$ belongs to $\boldsymbol{p}$. Let $n$ be the length of the longest path in $\mathcal{P}$, that is $\|\boldsymbol{p}\|_1 \leq n, \forall \boldsymbol{p} \in \mathcal{P}$. Given a time horizon $T \in \mathbb{N}$, at each (discrete) stage $t \in [T]$, a *learner* chooses a path $\tilde{\boldsymbol{p}}_t \in \mathcal{P}$. Then, a *loss vector* $\boldsymbol{\ell}_t \in [0,1]^E$ is secretly and adversarially chosen. Each element $\boldsymbol{\ell}_t(e)$ corresponds to the scalar loss embedded on the edge $e \in \mathcal{E}$. Note that we consider the *non-oblivious adversary*, i.e., $\boldsymbol{\ell}_t$ can be an arbitrary function of the learner's past actions $\tilde{\boldsymbol{p}}_s, \forall s \in [t-1]$, but not $\tilde{\boldsymbol{p}}_t$.[2] The learner's incurred loss is $L_t(\tilde{\boldsymbol{p}}_t) = (\tilde{\boldsymbol{p}}_t)^\top \boldsymbol{\ell}_t = \sum_{e \in \tilde{\boldsymbol{p}}_t} \boldsymbol{\ell}_t(e)$, i.e., the sum of the losses from the edges belonging to $\tilde{\boldsymbol{p}}_t$. The learner's feedback at stage $t$ after choosing $\tilde{\boldsymbol{p}}_t$ is presented as follows. First, she receives a *semi-bandit* feedback, that is, she observes the edges' losses $\boldsymbol{\ell}_t(e)$, for any $e$ belonging to the chosen path $\tilde{\boldsymbol{p}}_t$. Additionally, each edge $e \in \tilde{\boldsymbol{p}}_t$ may reveal the losses on several other edges. To represent these *side-observations* at time $t$, we consider a graph, denoted $G_t^O$, containing $E$ vertices. Each vertex $v_e$ of $G_t^O$ corresponds to an edge $e \in \mathcal{E}$ of the graph $G$. There exists a directed edge from a vertex $v_e$ to a vertex $v_{e'}$ in $G_t^O$ if, by observing the edge loss $\boldsymbol{\ell}_t(e)$, the learner can also deduce the edge loss $\boldsymbol{\ell}_t(e')$; we also denote this by $e \to e'$ and say that the edge $e$ reveals the edge $e'$. The objective of the learner is to minimize the cumulative *expected regret*, defined as $R_T := \mathbb{E}\left[\sum_{t \in [T]} L(\tilde{\boldsymbol{p}}_t)\right] - \min_{\boldsymbol{p}^* \in \mathcal{P}} \sum_{t \in [T]} L(\boldsymbol{p}^*)$.

Hereinafter, in places where there is no ambiguity, we use the term *path* to refer to a path in $\mathcal{P}$ and the term *observation graphs* to refer to $G_t^O$. In general, these observation graphs can depend on the decisions of both the learner and the adversary. On the other hand, all vertices in $G_t^O$ always have self-loops. In the case where none among $G_t^O, t \in [T]$ contains any other edge than these self-loops, no side-observation is allowed and the problem is reduced to the

---

[2]This setting is considered by most of the works in the non-stochastic/adversarial bandits literature, e.g., (Alon et al. 2013; Cesa-Bianchi and Lugosi 2012).

classical semi-bandit setting. If all $G_t^O, t \in [T]$ are complete graphs, SOPPP corresponds to the full-information PPPs. In this work, we focus on considering the *uninformed setting*, i.e., the learner observes $G_t^O$ only after making a decision at time $t$. On the other hand, we introduce two new notations:

$$\mathbb{O}_t(e) := \{\boldsymbol{p} \in \mathcal{P} : \exists e' \in \boldsymbol{p}, e' \to e\}, \forall e \in \mathcal{E},$$
$$\mathbb{O}_t(\boldsymbol{p}) := \{e \in \mathcal{E} : \exists e' \in \boldsymbol{p}, e' \to e\}, \forall \boldsymbol{p} \in \mathcal{P}.$$

Intuitively, $\mathbb{O}_t(e)$ is the set of all paths that, if chosen, reveal the loss on the edge $e$ and $\mathbb{O}_t(\boldsymbol{p})$ is the set of all edges whose losses are revealed if the path $\boldsymbol{p}$ is chosen. Trivially, $\boldsymbol{p} \in \mathbb{O}(e) \Leftrightarrow e \in \mathbb{O}(\boldsymbol{p})$. Moreover, due to the semi-bandit feedback, if $\boldsymbol{p}^* \ni e^*$, then $\boldsymbol{p}^* \in \mathbb{O}_t(e^*)$ and $e^* \in \mathbb{O}_t(\boldsymbol{p}^*)$. Apart from the results for general observation graphs, in this work, we additionally present several results under two particular assumptions, satisfied by some instances in practice (e.g., the CB and HS games), that provide more refined regret bounds compared to cases that were considered by (Kocák et al. 2014):

$(i)$ *symmetric* observation graphs where for each edge from $v_e$ to $v_{e'}$, there also exists an edge from $v_{e'}$ to $v_e$ (i.e., if $e \to e'$ then $e' \to e$); i.e., $G_t^O$ is an undirected graph;

$(ii)$ observation graphs under the following *assumption* $(A0)$ that requires that if two edges belong to a path in $G$, then they cannot simultaneously reveal the loss of another edge:
**Assumption $(A0)$:** *For any $e \in \mathcal{E}$, if $e' \to e$ and $e'' \to e$, then $\nexists \boldsymbol{p} \in \mathcal{P} : \boldsymbol{p} \ni e', \boldsymbol{p} \ni e''$.*

## 3 EXP3-OE - An Efficient Algorithm for the SOPPP

In this section, we present a new algorithm for SOPPP, called EXP3-OE (OE stands for Observable Edges), whose pseudo-code is given by Algorithm 1. The guarantees on the expected regret of EXP3-OE in SOPPP is analyzed in Section 3.2. Moreover, EXP3-OE always runs efficiently in polynomial time in terms of the number of edges of $G$; this is discussed in Section 3.1.

---

**Algorithm 1** EXP3-OE Algorithm for SOPPP.

1: **Input:** $T, \eta, \beta > 0$, graph $G$.
2: Initialize $w_1(e) := 1, \forall e \in \mathcal{E}$.
3: **for** $t = 1$ to $T$ **do**
4:     Loss vector $\boldsymbol{\ell}_t$ is chosen adversarially (unobserved).
5:     Use WP Algorithm (see Appendix A) to sample a path $\tilde{\boldsymbol{p}}_t$ according to $x_t(\tilde{\boldsymbol{p}}_t)$ (defined in (1)).
6:     Suffer the loss $L_t(\tilde{\boldsymbol{p}}_t) = \sum_{e \in \tilde{\boldsymbol{p}}_t} \boldsymbol{\ell}_t(e)$.
7:     Observation graph $G_t^O$ is generated and $\boldsymbol{\ell}_t(e)$, $\forall e \in \mathbb{O}_t(\tilde{\boldsymbol{p}}_t)$ are observed.
8:     $\hat{\boldsymbol{\ell}}_t(e) := \boldsymbol{\ell}_t(e) \mathbb{I}_{\{e \in \mathbb{O}_t(\tilde{\boldsymbol{p}}_t)\}} / (q_t(e) + \beta), \forall e \in \mathcal{E}$, where $q_t(e) := \sum_{\boldsymbol{p} \in \mathbb{O}_t(e)} x_t(\boldsymbol{p})$ is computed by Algorithm 2 (see Section 3.1).
9:     Update weights $w_{t+1}(e) := w_t(e) \cdot \exp(-\eta \hat{\boldsymbol{\ell}}_t(e))$.
10: **end for**

---

As an EXP3-type algorithm, EXP3-OE relies on the average weights sampling where at stage $t$ we update the

weight $w_t(e)$ on each edge $e$ by the exponential rule (line 9). For each path $\boldsymbol{p}$, we denote the path weight $w_t(\boldsymbol{p}) := \prod_{e \in \boldsymbol{p}} w_t(e)$ and define the following terms:

$$x_t(\boldsymbol{p}) := \frac{\prod\limits_{e \in \boldsymbol{p}} w_t(e)}{\sum\limits_{\boldsymbol{p}' \in \mathcal{P}} \prod\limits_{e' \in \boldsymbol{p}'} w_t(e')} = \frac{w_t(\boldsymbol{p})}{\sum\limits_{\boldsymbol{p}' \in \mathcal{P}} w_t(\boldsymbol{p}')}, \forall \boldsymbol{p} \in \mathcal{P}. \quad (1)$$

Line 5 of EXP3-OE involves a sub-algorithm, called the WPS algorithm, that samples a path $\boldsymbol{p} \in \mathcal{P}$ with probability $x_t(\boldsymbol{p})$ (the sampled path is then denoted by $\tilde{\boldsymbol{p}}_t$) from any input $\{w_t(e), e \in \mathcal{E}\}$ at each stage $t$. This algorithm is based on a classical technique called weight pushing (see e.g., (Takimoto and Warmuth 2003; György et al. 2007)). We discuss further details and present an explicit formulation of the WPS algorithm in Appendix A).

Compared to other instances of the EXP3-type algorithms, EXP3-OE has two major differences. First, at each stage $t$, the loss of each edge $e$ is estimated by $\hat{\ell}_t(e)$ (line 8) based on the term $q_t(e)$ and a parameter $\beta$. Intuitively, $q_t(e)$ is the probability that the loss on the edge $e$ is revealed from playing the chosen path at $t$. Second, the implicit exploration parameter $\beta$ added to the denominator allows us to "pretend to explore" in EXP3-OE without knowing the observation graph $G_t^O$ before making the decision at stage $t$ (the uninformed setting). Unlike the standard EXP3, the loss estimator used in EXP3-OE is *biased*, i.e., for any $e \in \mathcal{E}$,

$$\mathbb{E}_t\left[\hat{\ell}_t(e)\right] = \sum_{\tilde{\boldsymbol{p}} \in \mathcal{P}} x_t(\tilde{\boldsymbol{p}}) \frac{\ell_t(e)}{q_t(e) + \beta} \mathbb{I}_{\{e \in \mathbb{O}_t(\tilde{\boldsymbol{p}})\}}$$

$$= \sum_{\tilde{\boldsymbol{p}} \in \mathbb{O}_t(e)} x_t(\tilde{\boldsymbol{p}}) \frac{\ell_t(e)}{\sum\limits_{\boldsymbol{p} \in \mathbb{O}_t(e)} x_t(\boldsymbol{p}) + \beta} \leq \ell_t(e). \quad (2)$$

Here, $\mathbb{E}_t$ denotes the expectation w.r.t. the randomness of choosing a path at stage $t$. Second, unlike standard EXP3 algorithms that keep track and update on the weight of each path, the weight pushing technique is applied at line 5 (via the WPS algorithm) and line 8 (via Algorithm 2 in Section 3.1) where we work with edges weights instead of paths weights (recall that $E \ll P$).

### 3.1 Running Time Efficiency of the EXP3-OE Algorithm

In the WPS algorithm mentioned above, it is needed to compute the terms $H_t(s, u) := \sum_{\boldsymbol{p} \in \mathcal{P}_{s,u}} \prod_{e \in \boldsymbol{p}} w_t(e)$ and $H_t(u, d) := \sum_{\boldsymbol{p} \in \mathcal{P}_{u,d}} \prod_{e \in \boldsymbol{p}} w_t(e)$ for any vertex $u$ in $G$. Intuitively, $H_t(u, v)$ is the aggregate weight of all paths from vertex $u$ to vertex $v$ at stage $t$. These terms can be computed recursively in $\mathcal{O}(E)$ time based on dynamic programming. This computation is often referred to as weight pushing. Following the literature, we present in Appendix A an explicit algorithm that outputs $H_t(s, u), H_t(u, d), \forall u$ from any input $\{w_t(e), e \in \mathcal{E}\}$, called the WP algorithm. Then, a path in $G$ is sampled sequentially edge-by-edge based on these terms by the WPS algorithm. Importantly, the WP and WPS algorithms run efficiently in $\mathcal{O}(E)$ time.

The final non-trivial step to efficiently implement EXP3-OE is to compute $q_t(e)$ in line 8, i.e., the probability that

---

**Algorithm 2** Compute $q_t(e)$ of an edge $e$ at stage $t$.

1: **Input:** $e \in \mathbb{O}_t(\tilde{\boldsymbol{p}}_t)$, set $\mathfrak{R}_t(e)$ and $w_t(\bar{e}), \forall \bar{e} \in \mathcal{E}$.
2: Initialize $\bar{w}(\bar{e}) := w_t(\bar{e}), \forall \bar{e} \in \mathcal{E}$ and $q_t(e) := 0$.
3: Compute $H^*(s, d)$ by WP Algorithm (see Appendix A) with input $\{w_t(\bar{e}), \bar{e} \in \mathcal{E}\}$.
4: **for** $e' \in \mathfrak{R}_t(e)$ **do**
5:     Compute $H(s, u), H(u, d), \forall u \in \mathcal{V}$ by WP Algorithm with input $\{\bar{w}(\bar{e}), \forall \bar{e} \in \mathcal{E}\}$.
6:     $K(e') := H(s, u_{e'}) \cdot w(e') \cdot H(v_{e'}, d)$ where edge $e'$ goes from $u_{e'}$ to $v_{e'} \in C(u_{e'})$.
7:     $q_t(e) := q_t(e) + K(e')/H^*(s, d)$.
8:     Update $\bar{w}(e') = 0$.
9: **end for**
10: **Output:** $q_t(e)$.

---

an edge $e$ is revealed at stage $t$. Note that $q_t(e)$ is the sum of $|\mathbb{O}_t(e)| = \mathcal{O}(P)$ terms; therefore, a direct computation is inefficient while a naive application of the weight pushing technique can easily lead to errors. To compute $q_t(e)$, we propose Algorithm 2, a non-straightforward application of weight pushing, in which we consecutively consider all the edges $e' \in \mathfrak{R}_t(e) := \{e' \in \mathcal{E} : e' \to e\}$. Then, we take the sum of the terms $x_t(\boldsymbol{p})$ of the paths $\boldsymbol{p}$ going through $e'$ by the weight pushing technique while making sure that each of these terms $x_t(\boldsymbol{p})$ is included only once, even if $\boldsymbol{p}$ has more than one edge revealing $e$ (this is a non-trivial step). In Algorithm 2, we denote by $C(u)$ the set of the direct successors of any vertex $u \in \mathcal{V}$. We give a proof that Algorithm 2 outputs exactly $q_t(e)$ as defined in line 8 of Algorithm 1 in Appendix B. Algorithm 2 runs in $\mathcal{O}(|\mathfrak{R}_t(e)|E)$ time; therefore, line 8 of Algorithm 1 can be done in at most $\mathcal{O}(E^3)$ time.

In conclusion, EXP3-OE runs in at most $\mathcal{O}(E^3T)$ time, this guarantee works even for the worst-case scenario. For comparison, the FPL-IX algorithm runs in $\mathcal{O}(E|\mathcal{V}|^2T)$ time in expectation and in $\tilde{\mathcal{O}}(n^{1/2}E^{3/2}\ln(E/\delta)T^{3/2})$ time with a probability at least $1 - \delta$ for an arbitrary $\delta > 0$.[3] That is, FPL-IX might fail to terminate with a strictly positive probability[4] and it is not guaranteed to have efficient running time in all cases. Moreover, although this complexity bound of FPL-IX is slightly better in terms of $E$, the complexity bound of EXP3-OE improves that by a factor of $\sqrt{T}$. As is often the case in no-regret analysis, we consider the setting where T is significantly larger than other parameters of the problems; this is also consistent with the motivational applications of the CB and HS games presented in Section 1. Therefore, our contribution in improving the algorithm's running time in terms of $T$ is relevant.

### 3.2 Performance of the EXP3-OE Algorithm

In this section, we present an upper-bound of the expected regret achieved by the EXP3-OE algorithm in the SOPPP.

---

[3] If one runs FPL-IX with Dijkstra's algorithm as the optimization oracle and with parameters chosen by (Kocák et al. 2014)

[4] A stopping criterion for FPL-IX can be chosen to avoid this issue but it raises the question on how one chooses the criterion such that the regret guarantees hold.

For the sake of brevity, with $x_t(\boldsymbol{p})$ defined in (1), for any $t \in [T]$ and $e \in \mathcal{E}$, we denote:

$$r_t(e) := \sum_{\boldsymbol{p} \ni e} x_t(\boldsymbol{p}) \text{ and } Q_t := \sum_{e \in \mathcal{E}} r_t(e)/(q_t(e)+\beta).$$

Intuitively, $r_t(e)$ is the probability that the chosen path at stage $t$ contains an edge $e$ and $Q_t$ is the summation over all the edges of the ratio of this quantity and the probability that the loss of an edge is revealed (plus $\beta$). We can bound the expected regret with this key term $Q_t$.

**Theorem 3.1.** *The expected regret of the* EXP3-OE *algorithm in the* SOPPP *satisfies:*

$$R_T \le \ln(P)/\eta + \left[\beta + (n \cdot \eta)/2\right] \cdot \sum_{t \in [T]} Q_t. \quad (3)$$

A complete proof of Theorem 3.1 can be found in Appendix C and has an approach similar to (Alon et al. 2013; Cesa-Bianchi and Lugosi 2012) with several necessary adjustments to handle the new biased loss estimator in EXP3-OE. To see the relationship between the structure of the side-observations of the learner and the bound of the expected regret, we look for the upper-bounds of $Q_t$ in terms of the observation graphs' parameters. Let $\alpha_t$ be the independence number[5] of $G_t^O$, we have the following statement.

**Theorem 3.2.** *Let us define* $M := \lceil 2E^2/\beta \rceil$, $N_t := \ln\left(1 + \frac{M+E}{\alpha_t}\right)$ *and* $K_t := \ln\left(1 + \frac{nM+E}{\alpha_t}\right)$. *Upper-bounds of* $Q_t$ *in different cases of* $G_t^O$ *are given in the following table:*

|  | SATISFIES $(A0)$ | NOT SATISFIES $(A0)$ |
|---|---|---|
| SYMMETRIC | $\alpha_t$ | $n\alpha_t$ |
| NON-SYMMETRIC | $1 + 2\alpha_t N_t$ | $2n(1 + \alpha_t K_t)$ |

A proof of this theorem is given in Appendix E. The main idea of this proof is based on several graph theoretical lemmas that are extracted from (Alon et al. 2013; Kocák et al. 2014; Mannor and Shamir 2011). These lemmas establish the relationship between the independence number of a graph and the ratios of the weights on the graph's vertices that have similar forms to the key-term $Q_t$. The case where observation graphs are non-symmetric and do not satisfy assumption $(A0)$ is the most general setting. Moreover, as showed in Theorem 3.2, the bounds of $Q_t$ are improved if the observation graphs satisfy either the symmetry condition or assumption $(A0)$. Intuitively, given the same independence numbers, a symmetric observation graph gives the learner more information than a non-symmetric one; thus, it yields a better bound on $Q_t$ and the expected regret. On the other hand, assumption $(A0)$ is a technical assumption that allows the use of different techniques in the proofs to obtain better bounds. These cases have not been explicitly analyzed in the literature while they are satisfied by several practical situations, including the CB and HS games (see Section 4).

---

[5] The independence number of a directed graph is computed while ignoring the direction of the edges.

Finally, we give results on the upper-bounds of the expected regret, obtained by the EXP3-OE algorithm, presented as a corollary of Theorems 3.1 and 3.2.

**Corollary 3.3.** *In* SOPPP, *let* $\alpha$ *be an upper bound of* $\alpha_t, \forall t \in [T]$. *With appropriate choices of the parameters* $\eta$ *and* $\beta$, *the expected regret of the* EXP3-OE *algorithm is:*

$(i)$ $R_T \le \tilde{\mathcal{O}}(n\sqrt{T\alpha \ln(P)})$ *in the general cases.*

$(ii)$ $R_T \le \tilde{\mathcal{O}}(\sqrt{nT\alpha \ln(P)})$ *if assumption* $(A0)$ *is satisfied by the observation graphs* $G_t^O, \forall t \in [T]$.

A proof of Corollary 3.3 and the choices of the parameters $\beta$ and $\eta$ (these choices are non-trivial) yielding these results will be given in Appendix F. We can extract from this proof several more explicit results as follows: in the general case, $R_T \le \mathcal{O}\left(n\sqrt{T\alpha \ln(P)[1 + \ln(\alpha + \alpha \ln(\alpha) + E)]}\right)$ when the observations graphs are non-symmetric and $R_T \le (3/2)n\sqrt{T\alpha \ln(P)} + \sqrt{nT\alpha}$ if they are all symmetric; on the other hand, in cases that all the observation graphs satisfy $(A0)$, $R_T \le \mathcal{O}\left(\sqrt{nT\alpha \ln(P)[1 + 2\ln(1 + E)]}\right)$ if the observations graphs are non-symmetric and $R_T \le 2\sqrt{nT\alpha \ln(P)} + \sqrt{T\alpha}$ if they are all symmetric.
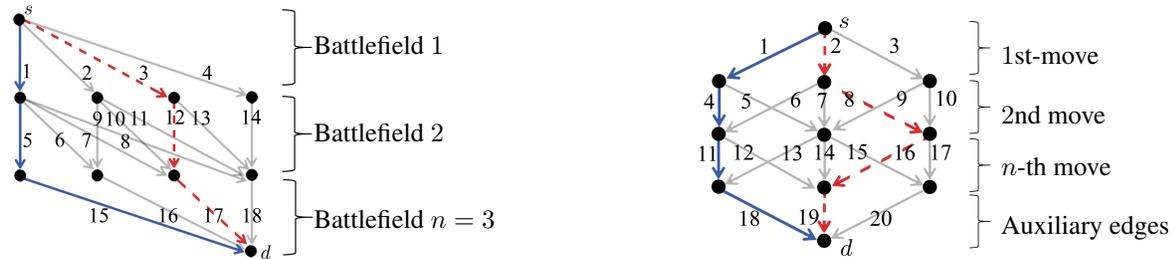
We note that a trivial upper-bound of $\alpha_t$ is the number of vertices of the graph $G_t^O$ which is $E$ (the number of edges in $G$). In general, the more connected $G_t^O$ is, the smaller $\alpha$ may be chosen; and thus the better upper-bound of the expected regret. In the (classical) semi-bandit setting, $\alpha_t = E, \forall t \in [T]$ and in the full-information setting, $\alpha_t = 1$, $\forall t \in [T]$. Finally, we also note that, if $P = \mathcal{O}(\exp(n))$ (this is typical in practice, including the CB and HS games), the bound in Corollary 3.3-$(i)$ matches in order with the bounds (ignoring the logarithmic factors) given by the FPL-IX algorithm (see (Kocák et al. 2014)). On the other hand, the form of the regret bound provided by the EXP3-IX algorithm (see (Kocák et al. 2014)) does not allow us to compare directly with the bound of EXP3-OE in the general SOPPP. EXP3-IX is only analyzed by (Kocák et al. 2014) when $n = 1$, i.e., $P = E$; in this case, we observe that the bound given by our EXP3-OE algorithm is better than that of EXP3-IX (by some multiplicative constants).

## 4 Colonel Blotto Games and Hide-and-Seek Games as SOPPP

Given the regret analysis of EXP3-OE in SOPPP, we now return to our main motivation, the Colonel Blotto and the Hide-and-Seek games, and discuss how to apply our findings to these games. To address this, we define formally the online version of the games and show how these problems can be formulated as SOPPP in Sections 4.1 and 4.2, then we demonstrate the benefit of using the EXP3-OE algorithm for learning in these games (Section 4.3).

### 4.1 Colonel Blotto Games as an SOPPP

**The online Colonel Blotto game** (the CB game). This is a game between a learner and an adversary over $n \ge 1$ battlefields within a time horizon $T > 0$. Each battlefield $i \in [n]$

(a) The graph $G_{3,3}$ corresponding to the CB game with $k=n=3$. E.g., the bold-blue path represents the strategy $(0,0,3)$ while the dash-red path represents the strategy $(2,0,1)$.

(b) The graph $G_{3,3,1}$ corresponding to the HS game with $k=n=3$ and $\kappa=1$. E.g., the blue-bold path represents the $(1,1,1)$ search and the red-dashed path represents the $(2,3,2)$ search.

Figure 1: Examples of the graphs corresponding to the CB game and the HS game.

has a value $\boldsymbol{b}_t(i) > 0$ (unknown to the learner)[6] at stage $t$ such that $\sum_{i=1}^n \boldsymbol{b}_t(i) = 1$. At stage $t$, the learner needs to distribute $k$ troops ($k \geq 1$ is fixed) towards the battlefields while the adversary simultaneously allocate hers; that is, the learner chooses a vector $\boldsymbol{z}_t$ in the strategy set $S_{k,n} := \{\boldsymbol{z} \in \mathbb{N}^n : \sum_{i=1}^n \boldsymbol{z}(i) = k\}$. At stage $t$ and battlefield $i \in [n]$, if the adversary's allocation is strictly larger than the learner's allocation $\boldsymbol{z}_t(i)$, the learner loses this battlefield and she suffers the loss $\boldsymbol{b}_t(i)$; if they have tie allocations, she suffers the loss $\boldsymbol{b}_t(i)/2$; otherwise, she wins and suffers no loss. At the end of stage $t$, the learner observes the loss from each battlefield (and which battlefield she wins, ties, or loses) but not the adversary's allocations. The learner's loss at each time is the sum of the losses from all the battlefields. The objective of the learner is to minimize her expected regret. Note that similar to SOPPP, we also consider the non-oblivious adversaries in the CB game.

While this problem can be formulated as a standard OComb, it is difficult to derive an efficient learning algorithm under that formulation, due to the learner's exponentially large set of strategies that she can choose from per stage. Instead, we show that by reformulating the problem as an SOPPP, we will be able to exploit the advantages of the Exp3-OE algorithm to solve it. To do so, first note that the learner can deduce several side-observations as follows: $(i)$ if she allocates $\boldsymbol{z}_t(i)$ troops to battlefield $i$ and wins, she knows that if she had allocated more than $\boldsymbol{z}_t(i)$ troops to $i$, she would also have won; $(ii)$ if she knows the allocations are tie at battlefield $i$, she knows exactly the adversary's allocation to this battlefield and deduce all the losses she might have suffered if she had allocated differently to battlefield $i$; $(iii)$ if she allocates $\boldsymbol{z}_t(i)$ troops to battlefield $i$ and loses, she knows that if she had allocated less than $\boldsymbol{z}_t(i)$ to battlefield $i$, she would also have lost.

Now, to cast the CB game as SOPPP, for each instance of the parameters $k$ and $n$, we create a DAG $G := G_{k,n}$ such that the strategy set $S_{k,n}$ has a one-to-one correspondence to the paths set $\mathcal{P}$ of $G_{k,n}$. Due to the lack of space, we only present here an example illustrating the graph of an instance of the CB game in Figure 1-(a) and we give the

formal definition of $G_{k,n}$ in Appendix G. The graph $G_{k,n}$ has $E = \mathcal{O}(k^2 n)$ edges and $P = |S_{k,n}| = \Omega\left(2^{\min\{n-1,k\}}\right)$ paths while the length of every path is $n$. Each edge in $G_{k,n}$ corresponds to allocating a certain amount of troops to a battlefield. Therefore, the CB game model is equivalent to a PPP where at each stage the learner chooses a path in $G_{k,n}$ and the loss on each edge is generated from the allocations of the adversary and the learner (corresponding to that edge) according to the rules of the game. At stage $t$, the (semi-bandit) feedback and the side-observations[7] deduced by the learner as described above infers an observation graph $G_t^O$. This formulation transforms any CB game into an SOPPP.

Note that since there are edges in $G_{m,n}$ that refer to the same allocation (e.g., the edges $5, 9, 12$, and $14$ in $G_{3,3}$ all refer to allocating $0$ troops to battlefield $2$), in the observation graphs, the vertices corresponding to these edges are always connected. Therefore, an upper bound of the independence number $\alpha_t$ of $G_t^O$ in the CB game is $\alpha_{\text{CB}} = n(k+1) = \mathcal{O}(nk)$. Moreover, we can verify that the observation graph $G_t^O$ of the CB game *satisfies assumption* $(A0)$ for any $t$ and it is *non-symmetric*.

### 4.2 Hide-and-Seek Games as an SOPPP

**The online Hide-and-Seek game** (the HS game). This is a repeated game (within the time horizon $T > 0$) between a hider and a seeker. In this work, we consider that the learner plays the role of the seeker and the hider is the adversary. There are $k$ locations, indexed from $1$ to $k$. At stage $t$, the learner sequentially chooses $n$ locations ($1 \leq n \leq k$), called an $n$-search, to seek for the hider, that is, she chooses $\boldsymbol{z}_t \in [k]^n$ (if $\boldsymbol{z}_t(i) = j$, we say that location $j$ is her $i$-th move). The hider maliciously assigns losses on all $k$ locations (intuitively, these losses can be the wasted time supervising a mismatch location or the probability that the hider does not hide there, etc.). In the HS game, the adversary is non-oblivious; moreover, in this work, we consider the following condition on how the hider/adversary assigns the losses on the locations:

---

[6]Knowledge on the battlefields' values is not assumed lest it limits the scope of application of our model (e.g., they are unknown in the radio resource allocation problem discussed in Section 1).

[7]E.g., in Figure 1-(a), if the learner chooses a path going through edge $10$ (corresponding to allocating $1$ troop to battlefield $2$) and wins (thus, the loss at edge $10$ is $0$), then she deduces that the losses on the edges $6, 7, 8, 10, 11$, and $13$ (corresponding to allocating at least $1$ troop to battlefield $2$) are all $0$.

(C1) *At stage $t$, the adversary secretly assigns a loss $\boldsymbol{b}_t(j)$ to each location $j \in [k]$ (unknown to the learner). These losses are fixed throughout the $n$-search of the learner.*

The learner's loss at stage $t$ is the sum of the losses from her chosen locations in the $n$-search at stage $t$, that is $\sum_{i\in[n],j\in[k]} \mathbb{I}_{\{\boldsymbol{z}_t(i)=j\}} \boldsymbol{b}_t(j)$. Moreover, often in practice the $n$-search of the learner needs to satisfy some constraints. In this work, as an example, we use the following constraint: $|\boldsymbol{z}_t(i) - \boldsymbol{z}_t(i+1)| \leq \kappa, \forall i \in [n]$ for a fixed $\kappa \in [0, k-1]$ (called the *coherence constraint*), i.e., the seeker cannot search too far away from her previously chosen location.[8] At the end of stage $t$, the learner only observes the losses from the locations she chose in her $n$-search, and her objective is to minimize her expected regret over $T$.

Similar to the case of the CB game, tackling the HS game as a standard OCOMB is computationally involved. As such, we follow the SOPPP formulation instead. To do this, we create a DAG $G := G_{k,n,\kappa}$ whose paths set has a one-to-one correspondence to the set containing all feasible $n$-search of the learner in the HS game with $k$ locations under $\kappa$-coherent constraint. Figure 1-(b) illustrates the corresponding graph of an instance of the HS game and we give a formal definition of $G_{k,n,\kappa}$ in Appendix G. The HS game is equivalent to the PPP where the learner chooses a path in $G_{k,n,\kappa}$ and edges' losses are generated by the adversary at each stage (note that to ensure all paths end at $d$, there are $n$ auxiliary edges in $G_{k,n,\kappa}$ that are always embedded with 0 losses). Note that there are $E = \mathcal{O}(k^2 n)$ edges and $P = \Omega(\kappa^{n-1})$ paths in $G_{k,n,\kappa}$. Moreover, knowing that the adversary follows condition $(C1)$, the learner can deduce the following side-observations: within a stage, the loss at each location remains the same no matter when it is chosen among the $n$-search, i.e., knowing the loss of choosing location $j$ as her $i$-th move, the learner knows all the loss if she chooses location $j$ as her $i'$-th move for any $i' \neq i$. The semi-bandit feedback and side-observations as described above generate the observation graphs $G_t^O$ (e.g., in Figure 1-(b), the edges 1, 4, 6, 11, and 13 represent that location 1 is chosen; thus, they mutually reveal each other). The independence number of $G_t^O$ is $\alpha_{\text{HS}} = k$ for any $t$. The observation graphs of the HS game are *symmetric* and *do not satisfy* $(A0)$. Finally, we consider a relaxation of condition $(C1)$:

(C2) *At stage $t$, the adversary assigns a loss $\boldsymbol{b}_t(j)$ on each location $j \in [k]$. For $i = 2, \ldots, n$, after the learner chooses, say location $j_i$, as her $i$-th move, the adversary can observe that and change the losses $\boldsymbol{b}_t(j)$ for any location that has not been searched before by the learner,[9] i.e., she can change the losses $\boldsymbol{b}_t(j), \forall j \notin \{j_1, \ldots, j_i\}$.*

By replacing condition $(C1)$ with condition $(C2)$, we can limit the side-observations of the learner: she can only de-

---

[8]Our results can be applied to HS games with other constraints, such as $\boldsymbol{z}_t(i) \leq \boldsymbol{z}_t(i+1), \forall i \in [n]$, i.e., she can only search forward; or, $\sum_{i\in[n]} \mathbb{I}_{\{\boldsymbol{z}_t(i)=k^*\}} \leq \kappa$, i.e., she cannot search a location $k^* \in [k]$ more than $\kappa$ times, etc.

[9]An interpretation is that by searching a location, the learner/seeker "discovers and secures" that location; therefore, the adversary/hider cannot change her assigned loss at that place.

duce that if $i_1 < i_2$, the edges in $G_{k,n,\kappa}$ representing choosing a location as the $i_1$-th move reveals the edges representing choosing that same location as the $i_2$-th move; but *not vice versa*. In this case, the observation graph $G_t^O$ is non-symmetric; however, its independence number is still $\alpha_{\text{HS}} = k$ as in the HS games with condition $(C1)$.

## 4.3 Performance of EXP3-OE in the Colonel Blotto and Hide-and-Seek Games

Having formulated the CB game and the HS game as SOPPPs, we can use the EXP3-OE algorithm in these games. From Section 3.1 and the specific graphs of the CB and HS game, we can deduce that EXP3-OE runs in at most $\mathcal{O}(k^6 n^3 T)$ time. We remark again that EXP3-OE's running time is linear in $T$ and efficient in all cases unlike when we run FPL-IX in the CB and HS games. Moreover, we can deduce the following result directly from Corollary 3.3:

**Corollary 4.1.** *The expected regret of the EXP3-OE algorithm satisfies:*

$(i)$ $R_T \leq \tilde{\mathcal{O}}(\sqrt{nT\alpha_{CB}\ln(P)}) = \tilde{\mathcal{O}}(\sqrt{Tn^3 k})$ *in the CB games with $k$ troops and $n$ battlefields.*

$(ii)$ $R_T \leq \tilde{\mathcal{O}}(n\sqrt{T\alpha_{HS}\ln(P)}) = \tilde{\mathcal{O}}(\sqrt{Tn^3 k})$ *in the HS games with $k$ locations and $n$-search.*

At a high-level, given the same scale on their inputs, the independence numbers of the observation graphs in HS games are smaller than in CB games (by a multiplicative factor of $n$). However, since assumption $(A0)$ is satisfied by the observation graphs of the CB games and not by the HS games, the expected regret bounds of the EXP3-OE algorithm in these games have the same order of magnitude. From Corollary 4.1, we note that in the CB games, the order of the regret bounds given by EXP3-OE is better than that of the FPL-IX algorithm (thanks to the fact that $(A0)$ is satisfied).[10] On the other hand, in the HS games with $(C1)$, the regret bounds of the EXP3-OE algorithm improves the bound of FPL-IX but they are still in the same order of the games' parameters (ignoring the logarithmic factors).[11] Note that the the regret bound of EXP3-OE in the HS game with Condition $(C1)$ (involving symmetric observation graphs) is slightly better than that in the HS game with Condition $(C2)$.

We also conducted several numerical experiments that compares the running time and the actual expected regret of EXP3-OE and FPL-IX in CB and HS games. The numerical results are in consistent with theoretical results in

---

[10]More explicitly, in the CB game, FPL-IX has a regret at most $\mathcal{O}\left(\ln(k^2 n^2 T)\sqrt{\ln(k^2 n)(k^2 n^4 + Cn^4 kT)}\right) = \tilde{\mathcal{O}}(\sqrt{Tn^4 k})$ (C is a constant indicated by (Kocák et al. 2014)) and EXP3-OE's regret bound is $\mathcal{O}\left(\sqrt{n^2 kT \cdot \min\{n-1, k\}[1+2\ln(1+k^2 n)]}\right)$ (if $n - 1 \leq k$, we can rewritten this bound as $\tilde{\mathcal{O}}(\sqrt{Tn^3 k})$).

[11]More explicitly, in HS games with $(C1)$, FPL-IX's regret is $\mathcal{O}\left(\ln(k^2 n^2 T)\sqrt{\ln(k^2 n)(k^2 n^4 + Cn^3 kT)}\right) = \tilde{\mathcal{O}}(Tn^3 k)$ and EXP3-OE's regret is $\mathcal{O}\left((3/2)\sqrt{n^3 kT\ln(k)} + \sqrt{nkT}\right) = \tilde{\mathcal{O}}(Tn^3 k)$ (similar results can be obtained for the HS games with $(C2)$).

this work. Our code for these experiments can be found at https://github.com/dongquan11/CB-HS.SOPPP.

Finally, we compare the regret guarantees given by our EXP3-OE algorithm and by the OSMD algorithm (see (Audibert, Bubeck, and Lugosi 2014))—the benchmark algorithm for OCOMB with semi-bandit feedback (although OSMD does not run efficiently in general): EXP3-OE is better than OSMD in CB games if $\mathcal{O}\left(n \cdot \ln\left(n^3 k^5 \sqrt{T}\right)\right) \leq k$; in HS games $(C1)$ if $\mathcal{O}(n \ln \kappa) \leq k$ and in the HS games with condition $(C2)$ if $n \cdot \ln \kappa \ln\left(n^4 k^5 \sqrt{T}\right) \leq \mathcal{O}(k)$. We give a proof of this statement in Appendix H. Intuitively, the regret guarantees of EXP3-OE is better than that of OSMD in the CB games where the learner's budget is sufficiently larger than the number of battlefields and in the HS games where the total number of locations is sufficiently larger than the number of moves that the learner can make in each stage.

## 5 Conclusion

In this work, we introduce the EXP3-OE algorithm for the path planning problem with semi-bandit feedback and side-observations. EXP3-OE is always efficiently implementable. Moreover, it matches the regret guarantees compared to that of the FPL-IX algorithm (EXP3-OE is better in some cases). We apply our findings to derive the first solutions to the online version of the Colonel Blotto and Hide-and-Seek games. This work also extends the scope of application of the PPP model in practice, even for large instances.

## References

Alon, N.; Cesa-Bianchi, N.; Gentile, C.; and Mansour, Y. 2013. From bandits to experts: A tale of domination and independence. In *the 27th Advances in Neural Information Processing Systems (NeurIPS)*, 1610–1618.

Alon, N.; Cesa-Bianchi, N.; Dekel, O.; and Koren, T. 2015. On-line learning with feedback graphs: Beyond bandits. In *the 28th Conference on Learning Theory (COLT)*, volume 40, 23–35.

Audibert, J.-Y.; Bubeck, S.; and Lugosi, G. 2014. Regret in online combinatorial optimization. *Mathematics of Operations Research* 39(1):31–45.

Auer, P.; Cesa-Bianchi, N.; Freund, Y.; and Schapire, R. E. 2002. The nonstochastic multiarmed bandit problem. *SIAM journal on computing* 32(1):48–77.

Behnezhad, S.; Dehghani, S.; Derakhshan, M.; Aghayi, M. T. H.; and Seddighin, S. 2017. Faster and simpler algorithm for optimal strategies of blotto game. In *the 31st AAAI Conference on Artificial Intelligence (AAAI)*, 369–375.

Bhattacharya, S.; Başar, T.; and Falcone, M. 2014. Surveillance for security as a pursuit-evasion game. In *the 5th International Conference on Decision and Game Theory for Security (GameSec)*, 370–379.

Blocki, J.; Christin, N.; Datta, A.; Procaccia, A. D.; and Sinha, A. 2013. Audit games. In *the 23rd International Joint Conference on Artificial Intelligence (IJCAI)*, 41–47.

Borel, E. 1921. La théorie du jeu et les équations intégrales à noyau symétrique. *Comptes rendus de l'Académie des Sciences* 173(1304-1308):58.

Bower, J. L., and Gilbert, C. G. 2005. *From resource allocation to strategy*. Oxford University Press.

Cesa-Bianchi, N., and Lugosi, G. 2012. Combinatorial bandits. *Journal of Computer and System Sciences* 78(5):1404–1422.

Chen, W.; Wang, Y.; and Yuan, Y. 2013. Combinatorial multi-armed bandit: General framework and applications. In *the 30th International Conference on Machine Learning (ICML)*, 151–159.

Chia, P. H. 2012. Colonel Blotto in web security. In *the 11th Workshop on Economics and Information Security, WEIS Rump Session*, 141–150.

Chien, S. F.; Zarakovitis, C. C.; Ni, Q.; and Xiao, P. 2019. Stochastic asymmetric blotto game approach for wireless resource allocation strategies. *IEEE Transactions on Wireless Communications*.

Freund, Y., and Schapire, R. E. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences* 55(1):119–139.

Gross, O., and Wagner, R. 1950. A continuous colonel blotto game. Technical report, RAND project air force Santa Monica CA.

György, A.; Linder, T.; Lugosi, G.; and Ottucsák, G. 2007. The on-line shortest path problem under partial monitoring. *Journal of Machine Learning Research* 8(Oct):2369–2403.

Hespanha, J. P.; Prandini, M.; and Sastry, S. 2000. Probabilistic pursuit-evasion games: A one-step nash approach. In *the 39th IEEE Conference on Decision and Control (CDC)*, 2272–2277.

Kalai, A., and Vempala, S. 2005. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences* 71(3):291–307.

Kocák, T.; Neu, G.; Valko, M.; and Munos, R. 2014. Efficient learning by implicit exploration in bandit problems with side observations. In *the 28th Advances in Neural Information Processing Systems (NeurIPS)*, 613–621.

Korzhyk, D.; Conitzer, V.; and Parr, R. 2010. Complexity of computing optimal stackelberg strategies in security resource allocation games. In *the 24th AAAI Conference on Artificial Intelligence (AAAI)*, 805–810.

Kovenock, D., and Roberson, B. 2012. Coalitional Colonel Blotto games with application to the economics of alliances. *Journal of Public Economic Theory* 14(4):653–676.

Mannor, S., and Shamir, O. 2011. From bandits to experts: On the value of side-observations. In *the 25th Advances in Neural Information Processing Systems (NeurIPS)*, 684–692.

Masucci, A. M., and Silva, A. 2014. Strategic resource allocation for competitive influence in social networks. In *the 52nd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 951–958.

Masucci, A. M., and Silva, A. 2015. Defensive resource allocation in social networks. In *the 54th IEEE Conference on Decision and Control (CDC)*, 2927–2932.

Navda, V.; Bohra, A.; Ganguly, S.; and Rubenstein, D. 2007. Using channel hopping to increase 802.11 resilience to jamming attacks. In *INFOCOM 2007. 26th IEEE International Conference on Computer Communications. IEEE*, 2526–2530. IEEE.

Roberson, B. 2006. The Colonel Blotto game. *Economic Theory* 29(1):2–24.

Sakaue, S.; Ishihata, M.; and Minato, S.-i. 2018. Efficient bandit combinatorial optimization algorithm with zero-suppressed binary decision diagrams. In *International Conference on Artificial Intelligence and Statistics*, 585–594.

Schwartz, G.; Loiseau, P.; and Sastry, S. S. 2014. The heterogeneous Colonel Blotto game. In *the 7th International Conference on Network Games, Control and Optimization (NetGCoop)*, 232–238.

Takimoto, E., and Warmuth, M. K. 2003. Path kernels and multiplicative updates. *Journal of Machine Learning Research* 4(Oct):773–818.

Vidal, R.; Shakernia, O.; Kim, H. J.; Shim, D. H.; and Sastry, S. 2002. Probabilistic pursuit-evasion games: theory, implementation, and experimental evaluation. *IEEE transactions on robotics and automation* 18(5):662–669.

Von Neumann, J. 1953. A certain zero-sum two-person game equivalent to the optimal assignment problem. *Contributions to the Theory of Games* 2:5–12.

Vu, D. Q.; Loiseau, P.; and Silva, A. 2018. Efficient computation of approximate equilibria in discrete colonel blotto games. In *the 27th International Joint Conference on Artificial Intelligence (IJCAI)*, 519–526.

Wang, Q., and Liu, M. 2016. Learning in hide-and-seek. *IEEE/ACM Transactions on Networking* 24(2):1279–1292.

Yavin, Y. 1987. Pursuit–evasion differential games with deception or interrupted observation. In *Pursuit-Evasion Differential Games*. Elsevier. 191–203.

Yucek, T., and Arslan, H. 2009. A survey of spectrum sensing algorithms for cognitive radio applications. *IEEE communications surveys & tutorials* 11(1):116–130.

Zhang, C.; Lesser, V.; and Shenoy, P. 2009. A multi-agent learning approach to online distributed resource allocation. In *the 21st International Joint Conference on Artificial Intelligence (IJCAI)*.

# Appendix
## A   Weight Pushing for Path Sampling

We re-visit some useful results in the literature. In this section, we consider a DAG $G$ with parameters as introduced in Section 2. For simplicity, we assume that each edge in $\mathcal{E}$ belongs to at least one path in $\mathcal{P}$. Let us respectively denote by $C(u)$ and $F(u)$ the set of the direct successors and the set of the direct predecessors of any vertex $u \in \mathcal{V}$. Moreover, let $e_{[u,v]}$ and $\mathcal{P}_{u,v}$ respectively denote the edge and the set of all paths from vertex $u$ to vertex $v$.

Let us consider a weight $w(e) > 0$ for each edge $e \in \mathcal{E}$. It is needed in the EXP3-OE algorithm to sample a path $\tilde{\boldsymbol{p}} \in \mathcal{P}$ with the probability:

$$x(\tilde{\boldsymbol{p}}) := \left[ \prod_{e \in \tilde{\boldsymbol{p}}} w(e) \right] \Big/ \left[ \sum_{\boldsymbol{p} \in \mathcal{P}} \prod_{e \in \boldsymbol{p}} w(e) \right]. \quad (4)$$

A direct computation and sampling from $x(\tilde{\boldsymbol{p}}), \forall \tilde{\boldsymbol{p}} \in \mathcal{P}$ takes $\mathcal{O}(P)$ time which is very inefficient. To efficiently sample the path, we first label the vertices set by $\mathcal{V} = \{s = u_0, u_1, \ldots, d = u_K\}$ such that if there exists an edge connecting $u_i$ to $u_j$ then $i < j$. We then define the following terms for each vertex $u \in \mathcal{V}$:

$$H(s, u) := \sum_{\boldsymbol{p} \in \mathcal{P}_{s,u}} \prod_{e \in \boldsymbol{p}} w(e) \text{ and } H(u, d) := \sum_{\boldsymbol{p} \in \mathcal{P}_{u,d}} \prod_{e \in \boldsymbol{p}} w(e).$$

Intuitively, $H(u, v)$ is the aggregate weight of all paths from vertex $u$ to vertex $v$ and $H(s, d)$ is exactly the denominator in (4). These terms $H(s, u)$ and $H(u, d), \forall u \in \mathcal{V}$ can be recursively computed by the WP algorithm (i.e., Algorithm 3) that runs in $\mathcal{O}(E)$ time, through dynamic programming. This is called *weight pushing* and it is used by (György et al. 2007; Sakaue, Ishihata, and Minato 2018; Takimoto and Warmuth 2003).

---

**Algorithm 3** WP Algorithm.

---

1: **Input:** Graph $G$, set of weights $\{w(e), e \in \mathcal{E}\}$.
2: Initialization $H(s, u_0) := H(u_K, d) := 1$.
3: **for** $k = 1$ **to** $K$ **do**
4:    $H(u_{K-k}, d) := \sum\limits_{v \in C(u_{K-k})} w(e_{[u_{K-k},v]}) H(v, d).$
5:    $H(s, u_k) := \sum\limits_{v \in F(u_k)} w(e_{[v,u_k]}) H(s, v).$
6: **end for**
7: **Output:** $H(s, u), H(u, d), \forall u \in \mathcal{V}$.

---

Based on the WP algorithm (i.e., Algorithm 3), we construct the WPS algorithm (i.e., Algorithm 4) that uses the weights $w(e), e \in \mathcal{E}$ as inputs and randomly outputs a path in $\mathcal{P}$. Intuitively, starting from the source vertex $s = u_0$, Algorithm 4 sequentially samples vertices by vertices based on the terms $H(u, v)$ computed by Algorithm 3. It is noteworthy that Algorithm 4 also runs in $\mathcal{O}(E)$ time and it is trivial to prove that the probability that a path $\boldsymbol{p}$ is sampled from Algorithm 4 matches exactly $d(\boldsymbol{p})$.

## B   Proof of Algorithm 2's Output

*Proof.* Fixing an edge $e \in \mathcal{E}$, we prove that when Algorithm 2 takes the edges weights $\{w_t(e), e \in \mathcal{E}\}$

as the input, it outputs exactly $q_t = \sum_{\boldsymbol{p} \in \mathbb{O}_t(e)} x_t(\boldsymbol{p})$. We note that if $e' \in \mathfrak{R}_t(e) := \{e' : e' \to e\}$, then $\{\boldsymbol{p} \in \mathcal{P} : \boldsymbol{p} \ni e'\} \subset \mathbb{O}_t(e)$.

We denote $|\mathfrak{R}_t(e)| = \rho_e$ and label the edges in the set $\mathfrak{R}_t(e)$ by $\{e_1, e_2, \ldots, e_{\rho_e}\}$. The for-loop in lines 4-8 of Algorithm 2 consecutively run with the edges in $R_t(e)$ as follows:

$(i)$ After the for-loop runs for $e_1$, we have $K(e_1) := \sum_{\boldsymbol{p} \ni e_1} \prod_{\bar{e} \in \boldsymbol{p}} \bar{w}(\bar{e}) = \sum_{\boldsymbol{p} \ni e_1} w_t(\boldsymbol{p})$; therefore, $q_t(e) = \sum_{\boldsymbol{p} \ni e_1} x_t(\boldsymbol{p})$ since $H^*(s, d) = \sum_{\boldsymbol{p} \in \mathcal{P}} w_t(\boldsymbol{p})$ computed from the original weights $w_t(\bar{e}), \bar{e} \in \mathcal{E}$. Due to line 8 that sets $\bar{w}(e_1) := 0$, henceforth in Algorithm 2, the weight $\bar{w}(\boldsymbol{p}) := \prod_{e \in \boldsymbol{p}} \bar{w}(e)$ of any path $\boldsymbol{p}$ that contains $e_1$ is set to 0.

$(ii)$ Let the for-loop run for $e_2$, we have $K(e_2) := \sum_{\boldsymbol{p} \ni e_2} \bar{w}(\boldsymbol{p}) = \sum_{\{\boldsymbol{p} \ni e_2\} \backslash \{\boldsymbol{p} \ni e_1\}} w_t(\boldsymbol{p})$ because any path $\boldsymbol{p} \ni e_1$ has the weight $\bar{w}(\boldsymbol{p}) = 0$. Therefore, $q_t(e) = \sum_{\boldsymbol{p} \ni e_1} x_t(\boldsymbol{p}) + \sum_{\{\boldsymbol{p} \ni e_2\} \backslash \{\boldsymbol{p} \ni e_1\}} x_t(\boldsymbol{p})$.

$(iii)$ Similarly, after the for-loop runs for $e_i$ (where $i \in \{3, \ldots, \rho_e\}$), we have:

$$q_t(e) = \sum_{k=1}^{i} \left( \sum_{\{\boldsymbol{p} \ni e_k\} \backslash \bigcup_{j<k} \{\boldsymbol{p} \ni e_j\}} x_t(\boldsymbol{p}) \right).$$

$(iv)$ Therefore, after the for-loop finishes running for every edge in $\mathfrak{R}_t(e)$; we have $q_t := \sum_{\boldsymbol{p} \in \mathbb{O}_t(e)} x_t(\boldsymbol{p})$ where each term $x_t(\boldsymbol{p})$ was only counted once even if $\boldsymbol{p}$ contains more than one edge that reveals the edge $e$.

$\square$

## C Proof of Theorem 3.1

**Theorem 3.1.** *The expected regret of the* EXP3-OE *algorithm in the* SOPPP *satisfies:*

$$R_T \leq \ln(P)/\eta + \left[\beta + (n \cdot \eta)/2\right] \cdot \sum_{t \in [T]} Q_t. \quad (3)$$

*Proof.* We first denote[12] $W_t := \sum_{\boldsymbol{p} \in \mathcal{P}} w_t(\boldsymbol{p}), \forall t \in [T]$. From line 9 of Algorithm 1, we trivially have:

$$w_{t+1}(\boldsymbol{p}) = w_t(\boldsymbol{p}) \cdot \exp(-\eta \hat{L}_t(\boldsymbol{p})), \forall \boldsymbol{p} \in \mathcal{P}, \forall t \in [T-1]. \quad (5)$$

---

[12]We recall that $w_t(\boldsymbol{p}) := \prod_{e \in \boldsymbol{p}} w_t(e)$.

---

**Algorithm 4** WPS Algorithm.
___
1: **Input:** Graph $G$, set of weights $\{w(e), e \in \mathcal{E}\}$.
2: $H(u, d), \forall u \in \mathcal{V}$ are computed by Algorithm 3.
3: Initialize $Q := \{s\}$, vertex $u := s$.
4: **while** $u \neq d$ **do**
5:    Sample a vertex $v$ from $\mathcal{C}(u)$ with probability $w(e_{[u,v]})H(v, d)/H(u, d)$.
6:    Add $v$ to the set $Q$ and update $u := v$.
7: **end while**
8: **Output:** $\tilde{\boldsymbol{p}} \in \mathcal{P}$ going through all the vertices in $Q$
___

We recall that $\hat{L}_t(\boldsymbol{p}) := \sum_{e \in \boldsymbol{p}} \hat{\ell}_t(e)$ and the notation $\mathbb{E}_t$ denoting the expectation w.r.t. to the randomness in choosing $\tilde{\boldsymbol{p}}_t$ in Algorithm 1 (i.e., w.r.t. the information up to time $t-1$). From (2), we have:

$$\mathbb{E}_t\left[\hat{L}_t(\boldsymbol{p})\right] \leq L_t(\boldsymbol{p}) := \sum_{e \in \boldsymbol{p}} \ell_t(e), \forall \boldsymbol{p} \in \mathcal{P}. \quad (6)$$

Under the condition that $0 < \eta$, we obtain:

$$\begin{aligned}
\frac{W_{t+1}}{W_t} &= \sum_{\boldsymbol{p} \in \mathcal{P}} \frac{w_{t+1}(\boldsymbol{p})}{W_t} \\
&= \sum_{\boldsymbol{p} \in \mathcal{P}} \frac{w_t(\boldsymbol{p}) \cdot \exp(-\eta \hat{L}_t(\boldsymbol{p}))}{W_t} \\
&= \sum_{\boldsymbol{p} \in \mathcal{P}} x_t(\boldsymbol{p}) \cdot \exp(-\eta \hat{L}_t(\boldsymbol{p}))) \\
&\leq \sum_{\boldsymbol{p} \in \mathcal{P}} \left[ x_t(\boldsymbol{p}) \left(1 - \eta \hat{L}_t(\boldsymbol{p}) + \frac{\eta^2}{2}(\hat{L}_t(\boldsymbol{p}))^2\right) \right] \\
&= 1 - \sum_{\boldsymbol{p} \in \mathcal{P}} \left[ x_t(\boldsymbol{p}) \left(\eta \hat{L}_t(\boldsymbol{p}) - \frac{\eta^2}{2}(\hat{L}_t(\boldsymbol{p}))^2\right) \right]. \quad (7)
\end{aligned}$$

Here, the second equality comes from (5) and the inequality comes from the fact that $\exp(-a) \leq 1 - a + a^2/2$ for $a := \eta \hat{L}_t(\boldsymbol{p}) \geq 0$. Now, we use the inequality $\ln(1 - y) \leq -y$, $\forall y < 1$ for $y := \sum_{\boldsymbol{p} \in \mathcal{P}} \left[ x_t(\boldsymbol{p}) \left(\eta \hat{L}_t(\boldsymbol{p}) - \frac{\eta^2}{2}(\hat{L}_t(\boldsymbol{p}))^2\right) \right]$,[13] then from (7), we obtain

$$\begin{aligned}
&\ln\left(\frac{W_{T+1}}{W_1}\right) \\
&= \sum_{t=1}^{T} \ln\left(\frac{W_{t+1}}{W_t}\right) \\
&\leq \sum_{t=1}^{T} \left(-\eta \sum_{\boldsymbol{p} \in \mathcal{P}} x_t(\boldsymbol{p}) \hat{L}_t(\boldsymbol{p}) + \frac{\eta^2}{2} \sum_{\boldsymbol{p} \in \mathcal{P}} x_t(\boldsymbol{p})(\hat{L}_t(\boldsymbol{p}))^2\right). \quad (8)
\end{aligned}$$

On the other hand, let us fix a path $\boldsymbol{p}^* \in \mathcal{P}$, then

$$\begin{aligned}
&\ln\left(\frac{W_{T+1}}{W_1}\right) \\
&\geq \ln\left(\frac{w_{T+1}(\boldsymbol{p}^*)}{W_1}\right) \\
&= \ln \frac{w_T(\boldsymbol{p}^*) \exp(-\eta \hat{L}_T(\boldsymbol{p}^*))}{P} \\
&= \ln \frac{w_{T-1}(\boldsymbol{p}^*) \exp(-\eta \hat{L}_T(\boldsymbol{p}^*) - \eta \hat{L}_{T-1}(\boldsymbol{p}^*))}{P} \\
&= -\eta \sum_{t=1}^{T} \hat{L}_t(\boldsymbol{p}^*) - \ln(P). \quad (9)
\end{aligned}$$

In the arguments leading to (9), we again use (5) and the fact that $w_1(\boldsymbol{p}) = 1, \forall \boldsymbol{p} \in \mathcal{P}$, including $w_1(\boldsymbol{p}^*)$. Therefore,

---

[13]We can easily check that $\eta \hat{L}_t(\boldsymbol{p}) - \eta^2 \hat{L}_t(\boldsymbol{p})^2/2 < 1$ for any $\eta > 0$ and thus, $\sum_{\boldsymbol{p} \in \mathcal{P}} \left[ x_t(\boldsymbol{p}) \left(\eta \hat{L}_t(\boldsymbol{p}) - \frac{\eta^2}{2}(\hat{L}_t(\boldsymbol{p}))^2\right) \right] < 1$.

combining (8) and (9) then dividing both sides by $\eta$, we have:

$$\sum_{t=1}^{T}\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\hat{L}_t(\boldsymbol{p})$$

$$\leq\frac{\ln(P)}{\eta}+\sum_{t=1}^{T}\hat{L}_t(\boldsymbol{p}^*)+\frac{\eta}{2}\sum_{t=1}^{T}\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})(\hat{L}_t(\boldsymbol{p}))^2. \quad (10)$$

Now, we take $\mathbb{E}_t$ on both sides of (10), then we apply (6) to obtain:

$$\sum_{t=1}^{T}\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\mathbb{E}_t[\hat{L}_t(\boldsymbol{p})]$$

$$\leq\frac{\ln(P)}{\eta}+\sum_{t=1}^{T}L_t(\boldsymbol{p}^*)+\frac{\eta}{2}\sum_{t=1}^{T}\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\mathbb{E}_t[\hat{L}_t(\boldsymbol{p})^2]. \quad (11)$$

Now, we look for a lower bound of $\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\mathbb{E}_t\left[\hat{L}_t(\boldsymbol{p})\right]$. For any fixed $\boldsymbol{p}\in\mathcal{P}$, we consider:

$$\mathbb{E}_t\left[\sum_{e\in\boldsymbol{p}}\hat{\boldsymbol{\ell}}_t(e)\right]=\sum_{\tilde{\boldsymbol{p}}\in\mathcal{P}}\left[x_t(\tilde{\boldsymbol{p}})\sum_{e\in\boldsymbol{p}}\left(\frac{\boldsymbol{\ell}_t(e)}{q_t(e)+\beta}\mathbb{I}_{\{e\in\mathbb{O}_t(\tilde{\boldsymbol{p}})\}}\right)\right]$$

$$=\sum_{e\in\boldsymbol{p}}\sum_{\tilde{\boldsymbol{p}}\in\mathbb{O}(e)}x_t(\tilde{\boldsymbol{p}})\frac{\boldsymbol{\ell}_t(e)}{q_t(e)+\beta}$$

$$=\sum_{e\in\boldsymbol{p}}\frac{q_t(e)\boldsymbol{\ell}_t(e)}{q_t(e)+\beta}. \quad (12)$$

Using (12) and recalling that $\boldsymbol{\ell}_t(e)\leq 1, \forall e\in\mathcal{E}$, we have:

$$\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\mathbb{E}_t\left[\hat{L}_t(\boldsymbol{p})\right]-\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})L_t(\boldsymbol{p})$$

$$=\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\sum_{e\in\boldsymbol{p}}\frac{q_t(e)\boldsymbol{\ell}_t(e)}{q_t(e)+\beta}-\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\sum_{e\in\boldsymbol{p}}\boldsymbol{\ell}_t(e)$$

$$=\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\sum_{e\in\boldsymbol{p}}\boldsymbol{\ell}_t(e)\left(\frac{q_t(e)}{q_t(e)+\beta}-1\right)$$

$$\geq-\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\sum_{e\in\boldsymbol{p}}\frac{\beta}{q_t(e)+\beta}$$

$$=-\beta\sum_{e\in\mathcal{E}}\frac{\sum_{\boldsymbol{p}\ni e}x_t(\boldsymbol{p})}{q_t(e)+\beta}$$

$$=-\beta Q_t. \quad (13)$$

Therefore, a lower bound of $\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\mathbb{E}_t\left[\hat{L}_t(\boldsymbol{p})\right]$ is $\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})L_t(\boldsymbol{p})-\beta Q_t$.

Now, we look for an upper bound of $\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\mathbb{E}_t\left[\hat{L}_t(\boldsymbol{p})^2\right]$. To do this, fix $\boldsymbol{p}\in\mathcal{P}$, we consider

$$\mathbb{E}_t\left[\hat{L}_t(\boldsymbol{p})^2\right]$$

$$=\mathbb{E}_t\left[\left(\sum_{e\in\boldsymbol{p}}\hat{\boldsymbol{\ell}}_t(e)\right)^2\right]$$

$$\leq n\cdot\mathbb{E}_t\left[\sum_{e\in\boldsymbol{p}}\hat{\boldsymbol{\ell}}_t(e)^2\right]$$

$$=n\cdot\sum_{\tilde{\boldsymbol{p}}\in\mathcal{P}}\left[x_t(\tilde{\boldsymbol{p}})\sum_{e\in\boldsymbol{p}}\left(\frac{\boldsymbol{\ell}_t(e)}{q_t(e)+\beta}\mathbb{I}_{\{e\in\mathbb{O}_t(\tilde{\boldsymbol{p}})\}}\right)^2\right]$$

$$\leq n\cdot\sum_{e\in\boldsymbol{p}}\sum_{\tilde{\boldsymbol{p}}\in\mathbb{O}_t(e)}x_t(\tilde{\boldsymbol{p}})\frac{1}{(q_t(e)+\beta)^2}$$

$$=n\cdot\sum_{e\in\boldsymbol{p}}q_t(e)\frac{1}{(q_t(e)+\beta)^2}$$

$$\leq n\cdot\sum_{e\in\boldsymbol{p}}\frac{1}{q_t(e)+\beta}. \quad (14)$$

The first inequality comes from applying Cauchy–Schwarz inequality. The second inequality comes from the fact that $\boldsymbol{\ell}_t(e)\leq 1$ and the last inequality comes from $q_t(e)\leq q_t(e)+\beta$ since $\beta > 0$.

Now, applying (14), we can bound

$$\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\mathbb{E}_t\left[\hat{L}_t(\boldsymbol{p})^2\right]\leq n\cdot\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})\sum_{e\in\boldsymbol{p}}\frac{1}{q_t(e)+\beta}$$

$$=n\cdot\sum_{e\in\mathcal{E}}\sum_{\boldsymbol{p}\ni e}x_t(\boldsymbol{p})\frac{1}{q_t(e)+\beta}$$

$$=n\cdot\sum_{e\in\mathcal{E}}\frac{r_t(e)}{q_t(e)+\beta}=n\cdot Q_t. \quad (15)$$

Here, we recall the notation $r_t(e)$ and $Q_t$ defined in Section 3.2. Replacing (13) and (15) into (11), we have that the following inequality holds for any $\boldsymbol{p}^*\in\mathcal{P}$.

$$\sum_{t=1}^{T}\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})L_t(\boldsymbol{p})-\sum_{t=1}^{T}\beta Q_t-\sum_{t=1}^{T}L_t(\boldsymbol{p}^*)$$

$$\leq\frac{\ln(P)}{\eta}+\frac{\eta}{2}\sum_{t=1}^{T}nQ_t.$$

Therefore, we conclude that

$$R_T=\sum_{t=1}^{T}\sum_{\boldsymbol{p}\in\mathcal{P}}x_t(\boldsymbol{p})L_t(\boldsymbol{p})-\sum_{t=1}^{T}L_t(\boldsymbol{p}^*)$$

$$\leq\frac{\ln(P)}{\eta}+\sum_{t=1}^{T}Q_t\left(n\frac{\eta}{2}+\beta\right).$$

$\square$

# D Lemmas on Graphs' Independence Numbers

In this section, we present some lemmas in graph theory that will be used in the next section to prove Theorem 3.2. Consider a graph $\tilde{G}$ whose vertices set and edges set are respectively denoted by $\tilde{\mathcal{V}}$ and $\tilde{\mathcal{E}}$. Let $\tilde{\alpha}$ be its independence number.

**Lemma D.1.** *Let $\tilde{G}$ be an directed graph and $I_v$ be the in-degree of the vertex $v \in \tilde{\mathcal{V}}$, then*

$$\sum_{v \in \tilde{\mathcal{V}}} [1/(1+I_v)] \leq 2\tilde{\alpha} \ln\left(1 + |\tilde{\mathcal{V}}|/\tilde{\alpha}\right).$$

A proof of this lemma can be found in Lemma 10 of (Alon et al. 2013).

**Lemma D.2.** *Let $\tilde{G}$ be a directed graph with self-loops and consider the numbers $k(v) \in [0,1], \forall v \in \tilde{\mathcal{V}}$ such that there exists $\gamma > 0$ and $\sum_{v \in \tilde{\mathcal{V}}} k(v) \leq \gamma$. For any $c > 0$, we have*

$$\sum_{v \in \tilde{\mathcal{V}}} \frac{k(v)}{\frac{1}{\gamma}\sum_{v' \to v} k(v') + c} \leq 2\gamma\tilde{\alpha} \ln\left(1 + \frac{\gamma\lceil|\tilde{\mathcal{V}}|^2/c\rceil + |\tilde{\mathcal{V}}|}{\tilde{\alpha}}\right) + 2\gamma.$$

A proof of this lemma can be found in Lemma 1 of (Kocák et al. 2014).

**Lemma D.3.** *Let $\tilde{G}$ be an undirected graph with self-loops and consider the numbers $k(v) \geq 0$, $v \in \mathcal{V}$. We have*

$$\sum_{v \in \tilde{\mathcal{V}}} \left[k(v)/\sum_{v' \to v} k(v')\right] \leq \tilde{\alpha}.$$

This lemma is extracted from Lemma 3 of (Mannor and Shamir 2011).

# E   Proof of Theorem 3.2

**Theorem 3.2.** *Let us define $M := \lceil 2E^2/\beta\rceil$, $N_t := \ln\left(1 + \frac{M+E}{\alpha_t}\right)$ and $K_t := \ln\left(1 + \frac{nM+E}{\alpha_t}\right)$. Upper-bounds of $Q_t$ in different cases of $G_t^O$ are given in the following table:*

|  | SATISFIES $(A0)$ | NOT SATISFIES $(A0)$ |
| --- | --- | --- |
| SYMMETRIC | $\alpha_t$ | $n\alpha_t$ |
| NON-SYMMETRIC | $1 + 2\alpha_t N_t$ | $2n(1 + \alpha_t K_t)$ |

**Case 1:** $G_t^O$ *does not satisfy assumption* $(A0)$. Fixing an edge $e$, due to the fact that $n$ is the length of the longest paths in $\mathcal{P}$, we have

$$nq_t(e) = n\sum_{\boldsymbol{p} \in \mathbb{O}_t(e)} x_t(\boldsymbol{p}) \geq \sum_{e' \to e}\sum_{\boldsymbol{p} \ni e'} x_t(\boldsymbol{p}) = \sum_{e' \to e} r_t(e')$$

$$\Rightarrow Q_t = \sum_{e \in \mathcal{E}} \frac{r_t(e)}{q_t(e) + \beta} \leq \sum_{e \in \mathcal{E}} \frac{r_t(e)}{\frac{1}{n}\sum_{e' \to e} r_t(e') + \beta}. \quad (16)$$

*Case 1.1:* If $G_t^O$ is a non-symmetric (i.e., directed) graph, we apply Lemma D.2 with $\gamma = n, c = \beta$ on the graph $\tilde{G} = G_t^O$ (whose vertices set $\tilde{\mathcal{V}}$ corresponds to the edges set $\mathcal{E}$ of $G$) and the numbers[14] $k(v_e) = r_t(e), \forall v_e \in \tilde{\mathcal{V}}$ (i.e., $\forall e \in \mathcal{E}$). We obtain the following inequality:

$$\sum_{e \in \mathcal{E}} \frac{r_t(e)}{\frac{1}{n}\sum_{e' \to e} r_t(e') + \beta} \leq 2n\alpha_t \ln\left(1 + \frac{n\lceil E^2/\beta\rceil + E}{\alpha_t}\right) + 2n.$$

---

[14] We verify that these numbers satisfy

$$\sum_{e \in \mathcal{E}} r_t(e) = \sum_{e \in \mathcal{E}}\sum_{\boldsymbol{p} \ni e} x_t(\boldsymbol{p}) = \sum_{\boldsymbol{p} \in \mathcal{P}}\sum_{e \in \boldsymbol{p}} x_t(\boldsymbol{p}) \leq \sum_{\boldsymbol{p} \in \mathcal{P}} nx_t(\boldsymbol{p}) = n.$$

*Case 1.2:* If $G_t^O$ is a symmetric (i.e. undirected) graph, we apply Lemma D.3 with the graph $\tilde{G} = G_t^O$ (whose vertices set $\tilde{\mathcal{V}}$ corresponds to the edges set $\mathcal{E}$ of the graph $G$) and the numbers $k(v_e) = r_t(e), \forall v_e \in \tilde{V}$ (i.e., $\forall e \in \mathcal{E}$) to obtain:

$$\sum_{e \in \mathcal{E}} \frac{r_t(e)}{\frac{1}{n}\sum_{e' \to e} r_t(e') + \beta} \leq n\sum_{e \in \mathcal{E}} \frac{r_t(e)}{\sum_{e' \to e} r_t(e')} \leq n\alpha_t.$$

**Case 2:** $G_t^O$ *satisfies assumption* $(A0)$. Under this assumption, $q_t(e) = \sum_{e' \to e} r_t(e')$ due to the definition of $\mathbb{O}_t(e)$. Therefore, $Q_t = \sum_{e \in \mathcal{E}} \left[r_t(e)/\left(\sum_{e' \to e} r_t(e') + \beta\right)\right]$.

*Case 2.1:* If $G_t^O$ is a non-symmetric (i.e., directed) graph. We consider a discretized version of $x_t(\boldsymbol{p})$ for any path $\boldsymbol{p} \in \mathcal{P}$ that is $\tilde{x}_t(\boldsymbol{p}) := k/M$ where $k$ is the unique integer such that $(k-1)/M \leq x_t(\boldsymbol{p}) \leq k/M$; thus, $\tilde{x}_t(\boldsymbol{p}) - 1/M \leq x_t(\boldsymbol{p}) \leq \tilde{x}_t(\boldsymbol{p})$.

Let us denote the discretized version of $r_t(e)$ by $\tilde{r}_t(e) := \sum_{\boldsymbol{p} \ni e} \tilde{x}_t(\boldsymbol{p})$. We deduce that $r_t(e) \leq \tilde{r}_t(e)$ and

$$\sum_{e' \to e} r_t(e) \geq \sum_{e' \to e}\left(\tilde{r}_t(e') - \frac{1}{M}\right) \geq \sum_{e' \to e} \tilde{r}_t(e') - \frac{E}{M}.$$

We obtain the bound:

$$Q_t = \sum_{e \in \mathcal{E}} \frac{r_t(e)}{\left(\sum_{e' \to e} r_t(e') + \beta\right)} \leq \sum_{e \in \mathcal{E}} \frac{\tilde{r}_t(e)}{\sum_{e' \to e} \tilde{r}_t(e') - E/M + \beta}. \quad (17)$$

We now consider the following inequality: If $a, b \geq 0$ and $a + b \geq B > A > 0$, then

$$\frac{a}{a + b - A} \leq \frac{a}{a + b} + \frac{A}{B - A}. \quad (18)$$

A proof of this inequality can be found in Lemma 12 of (Alon et al. 2013). Applying (18)[15] with $a = \tilde{r}_t(e)$, $b = \sum_{e' \to e, e' \neq e} \tilde{r}_t(e') + \beta$, $A = \frac{E}{M}$, and $B = \beta$ to (17),

$$Q_t \leq \sum_{e \in \mathcal{E}} \left(\frac{\tilde{r}_t(e)}{\sum_{e' \to e} \tilde{r}_t(e') + \beta} + \frac{E/M}{\beta - E/M}\right)$$

$$\leq \sum_{e \in \mathcal{E}} \frac{\tilde{r}_t(e)}{\sum_{e' \to e} \tilde{r}_t(e')} + 1. \quad (19)$$

The last inequality comes from the fact that $\frac{E}{M\beta - E} \leq \frac{E}{2E^2 - E} \leq \frac{1}{2E - 1} \leq \frac{1}{E}, \forall E \geq 1$.

Finally, we create an auxiliary graph $G_t^*$ such that:

(i) Corresponding to each edge $e$ in $G$ (i.e., each vertex $v_e$ in $G_t^O$), there is a clique, called $\mathbb{C}(e)$, in the auxiliary graph $G_t^*$ with $M\tilde{r}_t(e) \in \mathbb{N}$ vertices.

(ii) In each clique $\mathbb{C}(e)$ of $G_t^*$, all vertices are pairwise connected with length-two cycles. That is, for any $k, k' \in \mathbb{C}(e)$, there is an edge from $k$ to $k'$ and there is an edge from $k'$ to $k$ in $G_t^*$.

---

[15] Trivially, we can verify that $a + b \geq B$ and $B > A$ comes from the fact that $\beta \geq \beta\frac{1}{E} > \frac{E}{\lceil 2E^2/\beta\rceil}$.

$(iii)$ If $e \to e'$, i.e., there is an edge in $G_t^O$ connecting $v_e$ and $v_{e'}$; then in $G_t^*$, all vertices in the clique $\mathbb{C}(e)$ are connected to all vertices in $\mathbb{C}(e')$.

We observe that the independence number $\alpha_t$ of $G_t^O$ is equal to the independence number of $G_t^*$. Moreover, the in-degree of each vertex $k \in (e)$ in the graph $G_t^*$ is:

$$I_k^* = M\tilde{r}_t(e) - 1 + \sum_{e' \to e, e' \neq e} M\tilde{r}_t(e') = \sum_{e' \to e} M\tilde{r}_t(e') - 1. \tag{20}$$

Let us denote $V_t^*$ the set of all vertices in $G_t^*$, then we have:

$$\sum_{e \in \mathcal{E}} \frac{\tilde{r}_t(e)}{\sum_{e' \to e} \tilde{r}_t(e')} = \sum_{e \in \mathcal{E}} \frac{M\tilde{r}_t(e)}{\sum_{e' \to e} M\tilde{r}_t(e')} = \sum_{e \in \mathcal{E}} \sum_{k \in \mathbb{C}(e)} \frac{1}{I_k^* + 1}$$

$$= \sum_{k \in V_t^*} \frac{1}{\tilde{I}_k + 1} \leq 2\alpha_t \ln\left(1 + \frac{M+E}{\alpha_t}\right). \tag{21}$$

Here, the second equality comes from the fact that $|\mathbb{C}(e)| = M\tilde{r}_t(e)$ and (20). The inequality is obtained by applying Lemma D.1 to the graph $G_t^*$ and the fact that $|V_t^*| = \sum_{e \in \mathcal{E}} M\tilde{r}_t(e) \leq M \sum_{e \in \mathcal{E}} (r_t(e) + 1/M) \leq E + M$.

In conclusion, combining (19) and (21), we obtain the regret-upper bound as given in Theorem 3.2 for this case of the observation graph.

*Case 2.2:* Finally, if $G_t^O$ is a symmetric (i.e., undirected) graph, we again apply Lemma D.3 to the graph $\tilde{G} = G_t^O$ and the numbers $k(v_e) = r_t(e)$ to obtain that $Q_t \leq \sum_{e \in \mathcal{E}} \left[ r_t(e) / \sum_{e' \to e} r_t(e') \right] \leq \alpha_t$. $\qquad \square$

# F  Parameters Tuning for EXP3-OE: Proof of Corollary 3.3

In this section, we suggest a choice of $\beta$ and $\eta$ that guarantees the expected regret given in Corollary 3.3.

**Corollary 3.3.** *In* SOPPP*, let $\alpha$ be an upper bound of $\alpha_t, \forall t \in [T]$. With appropriate choices of the parameters $\eta$ and $\beta$, the expected regret of the* EXP3-OE *algorithm is:*

*(i)* $R_T \leq \tilde{\mathcal{O}}(n\sqrt{T\alpha\ln(P)})$ *in the general cases.*
*(ii)* $R_T \leq \tilde{\mathcal{O}}(\sqrt{nT\alpha\ln(P)})$ *if assumption $(A0)$ is satisfied by the observation graphs $G_t^O, \forall t \in [T]$.*

**Case 1: Non-symmetric (i.e. directed) observation graphs that do not satisfy assumption** $(A0)$. We find the parameters $\beta$ and $\eta$ such that $R_t \leq \tilde{\mathcal{O}}\left(n\sqrt{T\alpha}\right)$. We note that $\alpha_t \geq 1, \forall t \in [T]$; therefore, recalling that $\alpha$ is an upper bound of $\alpha_t$, from Theorem 3.1 and 3.2, we have:

$$R_T \leq \frac{\ln(P)}{\eta} + \sum_{t=1}^{T} \left(n\frac{\eta}{2} + \beta\right) 2n \left[1 + \alpha_t \ln\left(1 + \frac{nM+E}{\alpha_t}\right)\right]$$

$$\leq \frac{\ln(P)}{\eta} + T\left(n\frac{\eta}{2} + \beta\right) 2n \left[1 + \alpha \ln\left(\alpha + nM + E\right)\right]$$

$$= \frac{\ln(P)}{\eta} + \eta T n^2 \left[1 + \alpha \ln\left(\alpha + nM + E\right)\right]$$

$$+ 2\beta T n \left[1 + \alpha \ln\left(\alpha + nM + E\right)\right]. \tag{22}$$

Recalling that $M := \lceil 2E^2/\beta \rceil$, by choosing any

$$\beta \leq 1/\sqrt{Tn[1 + \alpha \ln(\alpha + n\lceil E^2/\beta \rceil + E)]}, \tag{23}$$

and $\eta = \sqrt{\ln(P)}/\sqrt{n^2 T [1 + \alpha \ln(\alpha + n\lceil E^2/\beta \rceil + E)]}$,

we obtain the bound:

$$R_T \leq 2n\sqrt{T\ln(P) \cdot [1 + \alpha \ln(\alpha + nM + E)]}$$
$$+ 2\sqrt{Tn[\alpha + \alpha \ln(\alpha + nM + E)]} \tag{24}$$
$$\leq \tilde{\mathcal{O}}\left(n\sqrt{T\alpha \ln(P)}\right).$$

In practice, as long as it satisfies (23), the larger $\beta$ is, the better upper-bounds that EXP3-OE gives. As an example that (23) always has at least one solution, we now prove that it holds with

$$\beta^* = \frac{-Tn^2E^2 + \sqrt{(Tn^2E^2)^2 + 4Tn(1 + \alpha \ln \alpha + E + n)}}{2Tn(1 + \alpha \ln \alpha + E + n)}. \tag{25}$$

Indeed, $\beta^* > 0$ and it satisfies:

$$\beta^{*2} \cdot Tn(1 + \alpha \ln \alpha + E + n) + \beta^* Tn^2 E^2 = 1.$$

$$\Rightarrow \beta^{*2} \cdot Tn(1 + \alpha \ln \alpha + E) + \beta^{*2} Tn^2 \left(\frac{E^2}{\beta^*} + 1\right) = 1$$

$$\Rightarrow \beta^{*2} \cdot Tn(1 + \alpha \ln \alpha + E) + \beta^{*2} Tn^2 \left\lceil \frac{E^2}{\beta^*} \right\rceil \leq 1$$

$$\Rightarrow \beta^* \leq \frac{1}{\sqrt{Tn(1 + \alpha \ln \alpha + E + nM)}}.$$

On the other hand, applying the inequality $\ln(1 + x) \leq x$, $\forall x \geq 0$, we have:

$$\frac{nM + E}{\alpha} \geq \ln\left(1 + \frac{nM + E}{\alpha}\right)$$

$$\Rightarrow \frac{nM + E}{\alpha} + \ln \alpha \geq \ln(\alpha + nM + E)$$

$$\Rightarrow nM + E + \alpha \ln \alpha + 1 \geq \alpha \ln(\alpha + nM + E) + 1$$

$$\Rightarrow \frac{1}{\sqrt{Tn(1 + \alpha \ln \alpha + nM + E)}} \leq \frac{1}{\sqrt{Tn(\alpha \ln(\alpha + nM + E) + 1)}}.$$

Therefore, $\beta^*$ satisfies (23). Finally, note that with the choice of $\beta = \beta^* = \Omega\left(nE^2/[1 + \alpha \ln \alpha + E + n]\right)$ as in (25), we have

$$M = \lceil 2E^2/\beta \rceil \leq \mathcal{O}([1 + \alpha \ln \alpha + E + n]/n).$$

Combining this with (24), we obtain the regret bound indicated in Section 3.2.

**Case 2: symmetric observation graphs that do not satisfy** $(A0)$. Trivially, we have that if $\beta := 1/\sqrt{n\alpha T}$ and $\eta = 2\sqrt{\ln(P)}/\sqrt{n^2\alpha T}$, then

$$R_T \leq \frac{\ln(P)}{\eta} + \left(n\frac{\eta}{2} + \beta\right) n\alpha T$$

$$= \frac{1}{2}n\sqrt{\alpha T \ln(P)} + n\sqrt{\alpha T \ln(P)} + \sqrt{n\alpha T} \tag{26}$$

$$\leq \tilde{\mathcal{O}}\left(n\sqrt{\alpha T \ln(P)}\right).$$

**Case 3: non-symmetric observation graphs** $G_t^O$ **satisfying assumption** $(A0)$, $\forall t$. We will prove that $R_T \leq \tilde{\mathcal{O}}\left(\sqrt{nT\alpha\ln(P)}\right)$ for any

$$\beta \leq 1/\sqrt{T\alpha[1+2\ln(1+\lceil E^2/\beta\rceil + E)]}, \quad (27)$$

$$\eta = 2\sqrt{\ln(P)}/\sqrt{Tn\alpha\left[1+2\ln(\alpha+M+E)\right]}. \quad (28)$$

Indeed, from Theorem 3.1 and 3.2, we have:

$$R_T \leq \frac{\ln(P)}{\eta} + \sum_{t=1}^{T}\left(n\frac{\eta}{2}+\beta\right)\left[1+2\alpha_t\ln\left(1+\frac{M+E}{\alpha_t}\right)\right]$$

$$\leq \frac{\ln(P)}{\eta} + \sum_{t=1}^{T}\left(n\frac{\eta}{2}+\beta\right)[\alpha+2\alpha\ln(1+M+E)]$$

$$= \frac{\ln(P)}{\eta} + \eta T\alpha\frac{n}{2}\left[1+2\ln(1+M+E)\right]$$
$$+ \beta T\alpha\left[1+2\ln(1+M+E)\right]. \quad (29)$$

We replace (27) and (28) into (29) and obtain:

$$R_T \leq \frac{3}{2}\sqrt{Tn\alpha\left[1+2\ln(1+M+E)\right]\cdot\ln(P)}$$

$$+ \sqrt{T\alpha\left[1+2\ln(1+M+E)\right]}. \quad (30)$$

$$\leq \tilde{\mathcal{O}}\left(\sqrt{n\alpha T\ln(P)}\right).$$

A choice for $\beta$ that satisfies (27) is

$$\beta^* := \frac{-T\alpha E^2 + \sqrt{(T\alpha E^2)^2 + T\alpha(3+2E)}}{T\alpha(3+2E)}. \quad (31)$$

Moreover, with this choice of $\beta^* = \Omega(E^2/(3+2E))$, we can deduce that $M := \lceil 2E^2/\beta^*\rceil \leq \mathcal{O}(3+2E)$. Combining this with (30), we obtain the regret bound indicated in Section 3.2.

**Case 4: all observation graphs are symmetric and satisfy** $(A0)$. From Theorem 3.1 and 3.2, we trivially have that if $\beta := 1/\sqrt{\alpha T}$ and $\eta = 2\sqrt{\ln(P)}/\sqrt{n\alpha T}$, then $R_T \leq 2\sqrt{n\alpha T\ln(P)} + \sqrt{\alpha T} \leq \tilde{\mathcal{O}}\left(\sqrt{n\alpha T\ln(P)}\right)$.

# G  Graphical Representation of the Games' Actions Sets

## G.1  The Actions Set of the Colonel Blotto Games

We give a description of the graph corresponding to the actions set of the learner in the CB game who distributes $k$ troops to $n$ battlefields.

**Definition G.1** (CB Graph). *The graph $G_{k,n}$ is a DAG that contains:*

*(i) $N := 2 + (k+1)(n-1)$ vertices arranged into $n+1$ layers. Layer $0$ and Layer $n$, each contains only one vertex, respectively labeled $s := (0,0)$–the source vertex and $d := (n,k)$–the destination vertex. Each Layer $i \in [n-1]$ contains $k+1$ vertices whose labels are ordered from left to right by $(i,0), (i,1), \ldots, (i,k)$.*

*(ii) There are directed edges from vertex $(0,0)$ to every vertex in Layer $1$ and edges from every vertex in Layer $n-1$ to vertex $(n,k)$. For $i \in \{1,2,\ldots,n-2\}$, there exists an edge connecting vertex $(i,j_1)$ (of Layer $i$) to vertex $(i+1,j_2)$ (of Layer $(i+1)$) if $k \geq j_2 \geq j_1 \geq 0$.*

Particularly, $G_{k,n}$ has $E = (k+1)\left[4+(n-2)(k+2)\right]/2 = \mathcal{O}(nk^2)$ edges and $P = \binom{n+k-1}{n-1} = \mathcal{O}(2^{\min\{n-1,k\}})$ paths going from vertex $s := (0,0)$ to vertex $d := (k,n)$. The edge connecting vertex $(i,j_1)$ to vertex $(i+1,j_2)$ for any $i \in \{0,1,\ldots,n-1\}$ represents allocating $(j_2-j_1)$ troops to battlefield $i+1$. Moreover, each path from $s$ to $d$ represents a strategy in $S_{k,n}$. This is formally stated in Proposition G.2.

**Proposition G.2.** *Given $k$ and $n$, there is a one-to-one mapping between the action set $S_{k,n}$ of the learner in the CB game (with $k$ troops and $n$ battlefields) and the set of all paths from vertex $s$ to vertex $d$ of the graph $G_{k,n}$.*

The proof of this proposition is trivial and can be intuitively seen in Figure 1-(a). We note that a similar graph is studied by (Behnezhad et al. 2017); however, it is used for a completely different purpose and it also contains more edges and paths than $G_{k,n}$ (that are not useful in this work).

## G.2  The Actions Set of the Hide-and-Seek game

We give a description of the graph corresponding to the actions set of the learner in the HS games with the $n$-search among $k$ locations and coherence constraints $|z_t(i) - z_t(i+1)| \leq \kappa, \forall i \in [n]$ for a fixed $\kappa \in [0, k-1]$.

**Definition G.3** (HS Graph). *The graph $G_{k,\kappa,n}$ is a DAG that contains:*

*(i) $N := 2 + kn$ vertices arranged into $n+2$ layers. Layer $0$ and Layer $(n+1)$, each contains only one vertex, respectively labeled $s$–the source vertex and $d$–the destination vertex. Each Layer $i \in \{1,\ldots,n\}$ contains $k$ vertices whose labels are ordered from left to right by $(i,1), (i,2), \ldots, (i,k)$.*

*(ii) There are directed edges from vertex $s$ to every vertex in Layer $1$ and edges from every vertex in Layer $n$ to vertex $d$. For $i \in \{1,2,\ldots,n-1\}$, there exists an edge connecting vertex $(i,j_1)$ to vertex $(i+1,j_2)$ if $|j_1 - j_2| \leq \kappa$.*

The graph $G_{k,\kappa,n}$ has $E = 2k+(n-1)\left[k+\kappa(2k-\kappa-1)\right] = \mathcal{O}(nk^2)$ edges and at least $\Omega(\kappa^{n-1})$ paths from $s$ to $d$. The edges ending at vertex $d$ are the auxiliary edges that are added just to guarantee that all paths end at $d$; these edges do not represent any intuitive quantity related to the game. For the remaining edges, any edge that ends at the vertex $(i,j)$ represents choosing the location $j$ as the $i$-th move. In other words, a path starting from $s$, passing by vertices $(1,j_1), (2,j_2), \ldots, (n,j_n)$ and ending at $d$ represents the $n$-search that chooses location $j_1$, then moves to location $j_2$, then moves to location $j_3$, and so on.

**Proposition G.4.** *Given $k$, $\kappa$ and $n$, there is a one-to-one mapping between the action set $S_{k,\kappa,n}$ of the learner in the HS game (with $n$-search among $k$ locations and coherence constraints with parameter $\kappa$) and the set of all paths from vertex $s$ to vertex $d$ of the graph $G_{k,\kappa,n}$.*

# H  EXP3-OE Algorithm and OSMD Algorithm in the CB and HS Games

$(i)$ As stated in Section 4, the observation graphs in the CB games are non-symmetric and they satisfy assumption $(A0)$. If we choose $\beta = \beta^*$ as in (31), then $\beta$ satisfies (27). Moreover, $\beta = \mathcal{O}(1/\sqrt{TnE})$; thus, $M = \mathcal{O}(E^2\sqrt{TnE})$. From (30), the expected regret of EXP3-OE in this case is bounded by $\mathcal{O}\sqrt{Tn(\alpha_{CB})\ln M \ln(P)}$ (recall that $\alpha_{CB} = kn$ is an upper bound of independence numbers of the observation graphs in the CB games). Therefore, to guarantee that this bound is better than the bound of the OSMD algorithm (that is $\sqrt{2TnE}$), the following inequality needs to hold:

$$\mathcal{O}\left(\alpha_{CB} \cdot \ln M \ln(P)\right) \leq E$$
$$\Rightarrow \mathcal{O}\left(nk \cdot \ln\left(E^2\sqrt{TnE}\right)\ln(2^n)\right) \leq nk^2$$
$$\Rightarrow \mathcal{O}\left(\ln\left(E^2\sqrt{TnE}\right)\ln(2^n)\right) \leq k$$
$$\Rightarrow \mathcal{O}\left(n\ln\left(n^3k^5\sqrt{T}\right)\right) \leq k.$$

$(ii)$ As stated in Section 4, the observation graphs in the HS games with condition $(C1)$ are symmetric and do not satisfy assumption $(A0)$. If we choose $\beta = 1/\sqrt{n\alpha T}$ then by (26), we have that $R_T$ is bounded by $\mathcal{O}\left(n\sqrt{\alpha_{HS}T\ln(P)}\right)$ (recall that $\alpha_{HS} = k$ is an upper bound of the independence numbers of the observation graphs in the HS games). Therefore, to guarantee that this bound is better than the bound of the OSMD algorithm in HS games, the following inequality needs to hold:

$$\mathcal{O}\left(\alpha_{HS} \cdot n\ln(P)\right) \leq E$$
$$\Rightarrow \mathcal{O}\left(k \cdot n\ln(P)\right) \leq nk^2$$
$$\Rightarrow \mathcal{O}\left(\ln(P)\right) \leq k$$
$$\Rightarrow \mathcal{O}\left(n\ln\kappa\right) \leq k.$$

$(iii)$ Finally, the observation graphs in the HS games with condition $(C2)$ are non-symmetric and do not satisfy assumption $(A0)$. Therefore, if we choose $\beta = \beta^*$ as in (25), then $\beta$ satisfies (23). In this case, $\beta = \mathcal{O}(1/\sqrt{TnE})$ and $M = \mathcal{O}(E^2\sqrt{TnE})$. Therefore, from (24), in this case, $R_T$ is bounded by $\mathcal{O}(n\sqrt{T\alpha_{HS}\ln\alpha_{HS}}\ln(nM))$. Therefore, to guarantee that this bound is better than the bound of OSMD (that is, $\sqrt{2TnE}$), the following inequality needs to hold:

$$\mathcal{O}\left(\alpha_{HS} \cdot n\ln nM \ln(P)\right) \leq E$$
$$\Rightarrow \mathcal{O}\left(nk\ln\left(\kappa^n\right)\ln(nE^2\sqrt{TnE})\right) \leq nk^2$$
$$\Rightarrow \mathcal{O}\left(n\ln\kappa\ln\left(n^4k^5\sqrt{T}\right)\right) \leq k.$$