

# A Two-Player Resource-Sharing Game with Asymmetric Information

Mevan Wijewardena, Michael J. Neely

## Abstract

This paper considers a two-player game where each player chooses a resource from a finite collection of options. Each resource brings a random reward. Both players have statistical information regarding the rewards of each resource. Additionally, there exists an information asymmetry where each player has knowledge of the reward realizations of different subsets of the resources. If both players choose the same resource, the reward is divided equally between them, whereas if they choose different resources, each player gains the full reward of the resource. We first implement the iterative best response algorithm to find an  $\epsilon$ -approximate Nash equilibrium for this game. This method of finding a Nash equilibrium may not be desirable when players do not trust each other and place no assumptions on the incentives of the opponent. To handle this case, we solve the problem of maximizing the worst-case expected utility of the first player. The solution leads to counter-intuitive insights in certain special cases. To solve the general version of the problem, we develop an efficient algorithmic solution that combines online convex optimization and the drift-plus penalty technique.

## Index Terms

Resource-sharing games, congestion games, potential games, worst-case utility maximization, drift-plus penalty method

## I. INTRODUCTION

We consider the following game with two players, A and B. There are  $n$  resources, each denoted by an integer between 1 and  $n$ . Each player selects a resource without knowledge about the other player's selection. The state of the game is described by the random vector

The authors are with the Electrical Engineering department at the University of Southern California.

This work was supported in part by one or more of: NSF CCF-1718477, NSF SpecEES 1824418.

$\mathbf{W} = (W_1, W_2, \dots, W_n)^\top$ , where  $W_k$  is the reward random variable of resource  $k$ . We assume  $W_k$  to be independent random variables for each  $1 \leq k \leq n$ , taking non-negative real values. If both players choose the same resource  $k$ , each gets a utility of  $W_k/2$ . If they choose different resources  $k, l$ , they receive utilities of  $W_k$  and  $W_l$ , respectively. It is assumed that the mean and the variance of  $W_k$  exist and are finite for each  $1 \leq k \leq n$ . Both players know the distribution of  $\mathbf{W}$ . Our formulation allows for an information asymmetry between the players. In particular,  $\{1, 2, \dots, n\}$  can be partitioned into four sets  $\{\mathcal{A}, \mathcal{B}, \mathcal{C}, \mathcal{AB}\}$  where only player A observes the realizations of  $W_k$  for  $k \in \mathcal{A}$ , only player B observes the realizations of  $W_k$  for  $k \in \mathcal{B}$ , no player observes the realizations of  $W_k$  for  $k \in \mathcal{C}$ , and both players observe the realizations of  $W_k$  for  $k \in \mathcal{AB}$ .

This game can be used to model different real-world scenarios where the agents have asymmetric information regarding the involved information structure. One classic example is the problem of Multiple-Access Control (MAC) in communication systems. Here, communication channels are accessed by multiple users, and the data rate of a channel is shared amongst the users who select it [1]. A channel can be shared using Time Division Multiple Access (TDMA) or Frequency Division Multiple Access (FDMA), where in TDMA, the channel is time-shared among the users [2], [3], whereas in FDMA, the channel is frequency-shared among the users [4]. In both cases, the total data rate supported by the channel can be considered the utility of the channel. The problem of information asymmetry arises since a user might have precise information regarding the total data rate offered by some channels but not others, and the known channels can be different for different users. On the other hand, the users in such a system cannot be trusted since the system may have malicious users (for instance, jammers) who focus on reducing the data rate available to genuine users.

Modified versions of this game apply to problems in economics. For instance, consider a firm that chooses a market to enter from a pool of market options. The chosen market may also be chosen by another firm. The reward of a market is the revenue it brings. Assume a simplified model where there exists a total revenue for each market, and the total revenue is divided equally among the firms entering the market. A reward known to all firms can be considered public information, while a reward known only to one firm is private information of that firm.

The game defined above can be viewed as a stochastic version of the class of games defined in [5], which are resource-sharing games, also known as congestion games. In resource-sharing

games, players compete for a finite collection of resources. In a single turn of the game, each player is allowed to select a subset of the collection of resources, where the allowed subsets make up the action space of the player. Each resource offers a reward to each player who selected the particular resource, where the reward offered depends on the number of players who selected it. The relationship between the reward offered to a player by a resource and the number of users selecting it is captured by the reward function of the resource. A player's utility is equal to the sum of the rewards offered by the resources in the subset selected by the player. In [5], it is established that the above game has a pure-strategy (deterministic) Nash equilibrium.

Although in the classical setting, these games ignore the stochastic nature of the rewards offered by the resources, the idea of resource-sharing games has been extended to different stochastic versions [6], [7]. Versions of the game with information asymmetry have been considered through the work of [8] in the context of Bayesian games, which considers the information design problem for resource-sharing with uncertainty. Similar Bayesian games have also been considered in [9], [10]. It should be noted that in general resource-sharing games, no conditions are placed on the reward functions of the resources. The special case where the reward functions are non-decreasing in the number of players selecting the resource is called a cost-sharing game [11]. These games are typically treated as games where a cost is minimized rather than a utility being maximized. In fair cost-sharing games, the cost of a resource is divided equally among the players selecting the resource. We consider a fair reward allocation model, where the reward of a resource is equally shared among the players selecting the resource. It should be noted that in this model, the players have opposite incentives compared to a fair-cost sharing model.

The work on resource-sharing games assumes that the players either cooperate or have the incentive to maximize a private or a social utility. It is interesting to consider a stochastic version of the game with asymmetric information between players who do not necessarily trust each other and who place no assumptions on the incentives of the opponents. In this context, the players have no signaling or external feedback and take actions based only on their personal knowledge of the reward realizations for a subset of the resource options. In this paper, we consider the above problem and limit our attention to the two-player singleton case, where each player can choose only one resource.

In the first part of the paper, we provide an iterative best response algorithm to find an

$\epsilon$ -approximate Nash equilibrium of the system. In the second part, we solve the problem of maximizing the worst-case expected utility of the first player. We solve the problem in two cases. The first case is when both players do not know the realizations of the reward random variables of any of the resources, in which case an explicit solution can be constructed. This case yields a counter-intuitive solution that provides insight into the problem. One such insight is that, while it is always optimal to choose from a subset of resources with the highest average rewards, within that subset, one chooses the higher-valued rewards with lower probability. For the second case, we solve the general version of the problem by developing an algorithm that leverages the online optimization technique [12], [13] and the drift-plus penalty method [14]. This algorithm generates a mixture of  $\mathcal{O}(1/\epsilon^2)$  pure strategies, which, when used in an equiprobable mixture, provides a utility within  $\epsilon$  of optimality on average. Below, we summarize our major contributions.

- We consider the problem of a two-player singleton stochastic resource-sharing game with asymmetric information. We first provide an iterative best response algorithm to find an  $\epsilon$ -approximate Nash equilibrium of the system. This equilibrium analysis uses potential game concepts.
- When the players do not trust each other and place no assumptions on the incentives of the opponent, we solve the problem of maximizing the worst-case expected utility of the first player using a novel algorithm that leverages techniques from online optimization and the drift-plus penalty methods. The algorithm developed can be used to solve the general unconstrained problem of finding the randomized decision  $\alpha \in \{1, 2, \dots, n\}$ , which maximizes  $\mathbb{E}\{h(\mathbf{x}; \Theta)\}$ , where  $\mathbf{x} \in \mathbb{R}^n$  with  $x_k = \mathbb{E}\{\Gamma_k \mathbb{1}_{\{\alpha=k\}}\}$ ,  $\Theta \in \mathbb{R}^m$  and  $\Gamma \in \mathbb{R}^n$  are non-negative random vectors with finite second moments, and  $h$  is a concave function such that  $\tilde{h}(\mathbf{x}) = \mathbb{E}\{h(\mathbf{x}; \Theta)\}$  is Lipschitz continuous, entry-wise non-decreasing and has bounded subgradients.
- We show our algorithm uses a mixture of only  $\mathcal{O}(1/\epsilon^2)$  pure strategies using a detailed analysis of the sample path of the related virtual queues (our preliminary work on this algorithm used a mixture of  $\mathcal{O}(1/\epsilon^3)$  pure strategies). Virtual queues are also used for constrained online convex optimization in [13], but our problem structure is different and requires a different and more involved treatment.

### A. Background on Resource-Sharing Games

The classical resource-sharing game defined in [5] is a tuple  $(\mathcal{M}, \mathcal{N}, \mathcal{T}, \mathbf{r})$ , where  $\mathcal{M}$  is a set of  $m$  players,  $\mathcal{N}$  is a set of  $n$  resources,  $\mathcal{T} = \mathcal{T}_1 \times \mathcal{T}_2 \times \dots \times \mathcal{T}_m$  where  $\mathcal{T}_j$  is the set of possible actions of player  $j$  (which is a subset of  $2^{\mathcal{N}}$ ), and  $\mathbf{r} = (r_1, r_2, \dots, r_n)$ , where  $r_i : \mathbb{N}_0 \rightarrow \mathbb{R}$  is the reward function of resource  $i$ . Here, we use the notation  $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$ . Each player has complete knowledge about the tuple  $(\mathcal{M}, \mathcal{N}, \mathcal{T}, \mathbf{r})$ , but they do not have knowledge of the actions chosen by other players. For an action profile  $\mathbf{a} = (a_1, a_2, \dots, a_m) \in \mathcal{T}$ , the count function  $\#$  is a function from  $\mathcal{N} \times \mathcal{T}$  to  $\mathbb{N}_0$  where  $\#(i, \mathbf{a}) = \sum_{k=1}^m \mathbb{1}_{\{i \in a_k\}}$ . In other words,  $\#(i, \mathbf{a})$  is the number of players choosing resource  $i$  under action profile  $\mathbf{a}$ . We call the quantity  $r_i(\#(i, \mathbf{a}))$  the *per-player reward* of resource  $i$  under action profile  $\mathbf{a}$ . The utility  $u_j$  of player  $j$  is a function from  $\mathcal{T}$  to  $\mathbb{R}$ , where  $u_j(\mathbf{a}) = \sum_{i=1}^n \mathbb{1}_{\{i \in a_j\}} r_i(\#(i, \mathbf{a}))$ . In other words,  $u_j(\mathbf{a})$  is the sum of the per-player rewards of the resources chosen by player  $j$  under action profile  $\mathbf{a}$ . Resource-sharing games fall under the general category of potential games [15]. Potential games are the class of games for which the change in reward of any player as a result of changing their strategy can be captured by the change in a global potential function.

Many game variations of the resource-sharing game have been studied [16]. Weighted resource-sharing games [17], games with player-dependent reward functions [18], and games with resources having preferences over players [19] are some of the extensions. Singleton games, where each player is allowed to choose only one resource, have also been explored explicitly in the literature [20], [21]. Some of the extensions of the classical resource-sharing game possess a pure Nash equilibrium in the singleton case. Two examples would be the games with player-specific reward functions for a resource [18] and the games with priorities where the resources have preferences over the players [19].

Resource-sharing games have been extended to several stochastic versions. For instance, ref. [6] considers the selfish routing problem with risk-averse players in a network with stochastic delays. The work of [7] considers two scenarios where, in the first scenario, each player participates in the game with a certain probability, and in the second scenario, the reward functions are stochastic. The problem of information asymmetry in resource-sharing games has been addressed through the work of [22], [8], [9], [10]. The work of [22] considers a network congestion game where the players have different information sets regarding the edges of the network. Further,

ref. [8] considers a scenario with a single random state  $\theta$ , which determines the reward functions. The realization of  $\theta$  is known to a game manager who strategically provides recommendations (signaling) to the players to minimize the social cost. An information asymmetry arises among the players in this case due to the actions of the game manager during private signaling, where the game manager provides player-specific recommendations.

Resource-sharing games appear in a variety of applications such as service chain composition [23], congestion control [24], network design [25], load balancing networks [26], [27], resource sharing in wireless networks [28], spectrum sharing [29], radio access selection [30], non-orthogonal multiple access [31], [32], network selection [33], [34], and migration of species [35].

Our formulation differs from the literature on resource-sharing games since we consider a scenario that is difficult to be analyzed using the standard equilibrium-based approaches. This is due to the fact that the players do not trust each other and place no assumptions on the incentives of the opponents, and they take action in the absence of a signaling mechanism or external feedback by just using their knowledge of the reward random variables. This motivates our formulation as a one-shot problem tackled using worst-case expected utility maximization.

### B. Notation

We use calligraphic letters to denote sets. Vectors and matrices are denoted by boldface characters. For integers  $n$  and  $m$ , we denote by  $[n : m]$  the inclusive set of integers between  $n$  and  $m$ . Given a vector  $\mathbf{w} \in \mathbb{R}^m$ ,  $w_k$  is used to denote the  $k$ -th element of  $\mathbf{w}$ ;  $\mathbf{w}_{k:l}$  for  $l \geq k$  represents the  $l - k + 1$  dimensional sub-vector  $(w_k, w_{k+1}, \dots, w_l)^\top$  of  $\mathbf{w}$ ; for a subset  $\mathcal{S}$  of integers from 1 to  $n$   $\{w_k; k \in \mathcal{S}\}$  represents the sub-vector of  $\mathbf{w}$  with index in  $\mathcal{S}$ . For  $\mathbf{z} \in \mathbb{R}^m$ , we use  $\|\mathbf{z}\|_2$ , and  $\|\mathbf{z}\|_\infty$  to denote the standard Euclidean norm (L2 norm), and the supremum norm ( $\max\{|z_1|, |z_2|, \dots, |z_m|\}$ ) of  $\mathbf{z}$ , respectively. For a function  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ , and  $\mathbf{z} \in \mathbb{R}^m$ , we use  $f'(\mathbf{z}) = (f'_1(\mathbf{z}), f'_2(\mathbf{z}), \dots, f'_m(\mathbf{z}))$  to denote a subgradient of  $f$  at  $\mathbf{z}$ .

## II. MATERIALS AND METHODS

The code used for the simulations is implemented using Python programming language in the notebook <https://rb.gy/wvt33>.

### III. FORMULATION

Denote  $\mathbf{X} = \{W_k; k \in \mathcal{A}\}$ ,  $\mathbf{Y} = \{W_k; k \in \mathcal{B}\}$ ,  $\mathbf{Z} = \{W_k; k \in \mathcal{AB}\}$ , and  $\mathbf{V} = \{W_k; k \in \mathcal{C}\}$ . Recall that  $\mathbf{X}$  is known only to player A,  $\mathbf{Y}$  is known only to player B,  $\mathbf{Z}$  is known to both players, and  $\mathbf{V}$  is known to neither. Let us define  $\mathcal{A}^c = [1 : n] \setminus \mathcal{A}$ , and  $\mathcal{B}^c = [1 : n] \setminus \mathcal{B}$ . Let  $|\mathcal{A}| = a$ ,  $|\mathcal{B}| = b$ ,  $|\mathcal{C}| = c$  and  $|\mathcal{AB}| = d$ . Therefore,  $a + b + c + d = n$ . Without loss of generality, we assume  $\mathcal{A} = [1 : a]$ ,  $\mathcal{B} = [a + 1 : a + b]$ ,  $\mathcal{C} = [a + b + 1 : a + b + c]$ , and  $\mathcal{AB} = [a + b + c + 1 : n]$ .

Let  $R^C(g^A, g^B)$  be the random variable representing the utility of player  $C \in \{A, B\}$ , given that player A uses strategy  $g^A$ , and player B uses strategy  $g^B$ . General strategies for players A and B can be represented by the Borel-measurable functions,

$$g^A : [0, 1) \times \mathbb{R}_{\geq 0}^{a+d} \rightarrow [1 : n], \quad (1)$$

$$g^B : [0, 1) \times \mathbb{R}_{\geq 0}^{b+d} \rightarrow [1 : n], \quad (2)$$

where

$$\alpha^A = g^A(U^A, \mathbf{X}, \mathbf{Z}), \quad (3)$$

$$\alpha^B = g^B(U^B, \mathbf{Y}, \mathbf{Z}), \quad (4)$$

are the resources chosen by players A and B, respectively. Here,  $U^A$  and  $U^B$  are independent randomization variables uniformly distributed in  $[0, 1)$  and independent of  $\mathbf{W}$ . A pure strategy for player A is a function  $g^A$  that does not depend on  $U^A$ , whereas a mixed strategy is a function  $g^A$  that depends on  $U^A$ . Hence, we drop the randomization variable when depicting a pure strategy. Pure strategies and mixed strategies for player B are defined similarly. Let  $\mathcal{S}^A$  and  $\mathcal{S}^B$  denote the sets of all possible strategies for players A and B, respectively.

It turns out that our analysis is simplified when  $\mathbf{Z}$  is fixed. Fixing  $\mathbf{Z}$  does not affect the symmetry between players A and B since  $\mathbf{Z}$  is observed by both players A and B. Hereafter, we conduct the analysis by considering all quantities conditioned on  $\mathbf{Z}$ .

Define

$$\begin{aligned} p_k^A &= \mathbb{E}\{\mathbb{1}_{\{\alpha^A=k\}} | \mathbf{Z}\} \text{ for } 1 \leq k \leq n, \\ q_k^A &= \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} | \mathbf{Z}\} \text{ for } k \in \mathcal{A}, \end{aligned} \quad (5)$$

and,

$$p_k^B = \mathbb{E}\{\mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \text{ for } 1 \leq k \leq n,$$

$$q_k^B = \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \text{ for } k \in \mathcal{B}. \quad (6)$$

Note that  $p_k^A$  and  $p_k^B$  are the conditional probabilities of players A and B choosing  $k$  given  $\mathbf{Z}$ . Define vectors  $\mathbf{p}^A = \{p_k^A; 1 \leq k \leq n\}$ ,  $\mathbf{q}^A = \{q_k^A; k \in \mathcal{A}\}$ ,  $\mathbf{p}^B = \{p_k^B; 1 \leq k \leq n\}$ , and  $\mathbf{q}^B = \{q_k^B; k \in \mathcal{B}\}$ . For  $1 \leq k \leq n$ , define  $E_k = \mathbb{E}\{W_k | \mathbf{Z}\}$ . Hence, we have

$$E_k = \begin{cases} W_k & \text{if } k \in \mathcal{AB}, \\ \mathbb{E}\{W_k\} & \text{otherwise,} \end{cases} \quad (7)$$

which uses the independence of  $W_k$  and  $\mathbf{Z}$  when  $k \notin \mathcal{AB}$ .

Note that the utility achieved by player A given the strategies  $g^A$  and  $g^B$  can be written as

$$R^A(g^A, g^B) = \sum_{k=1}^n W_k \left( \mathbb{1}_{\{\alpha^A=k\}} - \frac{1}{2} \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} \right). \quad (8)$$

Given the strategies  $g^A$  and  $g^B$ , we provide an expression for the expected utility of player A given  $\mathbf{Z}$ , where the expectation is over the random variables  $\mathbf{X}, \mathbf{Y}, \mathbf{V}$ , and the possibly random actions  $\alpha^A$  and  $\alpha^B$ . Taking expectations of (8) gives,

$$\begin{aligned} \mathbb{E}\{R^A(g^A, g^B) | \mathbf{Z}\} &= \sum_{k=1}^n \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} | \mathbf{Z}\} - \frac{1}{2} \sum_{k=1}^n \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \\ &= \sum_{k \in \mathcal{A}} \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} | \mathbf{Z}\} + \sum_{k \in \mathcal{A}^c} \mathbb{E}\{W_k | \mathbf{Z}\} \mathbb{E}\{\mathbb{1}_{\{\alpha^A=k\}} | \mathbf{Z}\} - \frac{1}{2} \sum_{k=1}^n \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \\ &= \sum_{k \in \mathcal{A}} q_k^A + \sum_{k \in \mathcal{A}^c} E_k p_k^A - \frac{1}{2} \sum_{k=1}^n \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\}. \end{aligned} \quad (9)$$

Note that given  $\mathbf{Z}$ , the random variables  $\alpha^A$  and  $\alpha^B$  are independent. Hence, we can split the last term (9) as follows,

$$\begin{aligned} \sum_{k=1}^n \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} &= \sum_{k \in \mathcal{A}} \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} + \sum_{k \in \mathcal{B}} \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \\ &+ \sum_{k \in \mathcal{C} \cup \mathcal{AB}} \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} = \sum_{k \in \mathcal{A}} \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} | \mathbf{Z}\} \mathbb{E}\{\mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \\ &+ \sum_{k \in \mathcal{B}} \mathbb{E}\{\mathbb{1}_{\{\alpha^A=k\}} | \mathbf{Z}\} \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} + \sum_{k \in \mathcal{C} \cup \mathcal{AB}} E_k \mathbb{E}\{\mathbb{1}_{\{\alpha^A=k\}} | \mathbf{Z}\} \mathbb{E}\{\mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \\ &= \sum_{k \in \mathcal{A}} q_k^A p_k^B + \sum_{k \in \mathcal{B}} p_k^A q_k^B + \sum_{k \in \mathcal{C} \cup \mathcal{AB}} E_k p_k^A p_k^B. \end{aligned} \quad (10)$$



#### IV. COMPUTING THE $\epsilon$ -APPROXIMATE NASH EQUILIBRIUM

This section focuses on finding an  $\epsilon$ -approximate Nash equilibrium of the game. Fix  $\epsilon > 0$ . A strategy pair  $(g^A, g^B)$  is defined as an  $\epsilon$ -approximate Nash equilibrium if neither player can improve its expected reward by more than  $\epsilon$  if it changes its strategy (while holding the strategy of the other player fixed).

Combining (10) with (9), we have that

$$\mathbb{E}\{R^A(g^A, g^B)|\mathbf{Z}\} = \sum_{k \in \mathcal{A}} q_k^A + \sum_{k \in \mathcal{A}^c} E_k p_k^A - \frac{1}{2} \left( \sum_{k \in \mathcal{A}} q_k^A p_k^B + \sum_{k \in \mathcal{B}} p_k^A q_k^B + \sum_{k \in \mathcal{C} \cup \mathcal{A}\mathcal{B}} E_k p_k^A p_k^B \right). \quad (11)$$

Similarly, for player B, we have

$$\mathbb{E}\{R^B(g^A, g^B)|\mathbf{Z}\} = \sum_{k \in \mathcal{B}} q_k^B + \sum_{k \in \mathcal{B}^c} E_k p_k^B - \frac{1}{2} \left( \sum_{k \in \mathcal{A}} q_k^A p_k^B + \sum_{k \in \mathcal{B}} p_k^A q_k^B + \sum_{k \in \mathcal{C} \cup \mathcal{A}\mathcal{B}} E_k p_k^A p_k^B \right). \quad (12)$$

First, we focus on finding the best response for players A and B, given the other player's strategy is fixed.

*Lemma 1:* The best response for players A and B are given by  $\alpha^A = \arg \max_{1 \leq k \leq n} A_k$ , and  $\alpha^B = \arg \max_{1 \leq k \leq n} B_k$ , where  $A_k$  and  $B_k$  are given by,

$$A_k = \begin{cases} W_k \left(1 - \frac{p_k^B}{2}\right) & \text{if } k \in \mathcal{A}, \\ E_k - \frac{q_k^B}{2} & \text{if } k \in \mathcal{B}, \\ E_k \left(1 - \frac{p_k^B}{2}\right) & \text{if } k \in \mathcal{C} \cup \mathcal{A}\mathcal{B}, \end{cases} \quad B_k = \begin{cases} E_k - \frac{q_k^A}{2} & \text{if } k \in \mathcal{A}, \\ W_k \left(1 - \frac{p_k^A}{2}\right) & \text{if } k \in \mathcal{B}, \\ E_k \left(1 - \frac{p_k^A}{2}\right) & \text{if } k \in \mathcal{C} \cup \mathcal{A}\mathcal{B}. \end{cases} \quad (13)$$

*Proof:* We find the best response for A, and the best response for B follows similarly. Notice that we can rearrange (11) as,

$$\begin{aligned} \mathbb{E}\{R^A(g^A, g^B)|\mathbf{Z}\} &= \sum_{k \in \mathcal{A}} q_k^A \left(1 - \frac{p_k^B}{2}\right) + \sum_{k \in \mathcal{B}} p_k^A \left(E_k - \frac{q_k^B}{2}\right) + \sum_{k \in \mathcal{C} \cup \mathcal{A}\mathcal{B}} p_k^A E_k \left(1 - \frac{p_k^B}{2}\right) \\ &= \sum_{k \in \mathcal{A}} \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}}|\mathbf{Z}\} \left(1 - \frac{p_k^B}{2}\right) + \sum_{k \in \mathcal{B}} \mathbb{E}\{\mathbb{1}_{\{\alpha^A=k\}}|\mathbf{Z}\} \left(E_k - \frac{q_k^B}{2}\right) \\ &\quad + \sum_{k \in \mathcal{C} \cup \mathcal{A}\mathcal{B}} E_k \mathbb{E}\{\mathbb{1}_{\{\alpha^A=k\}}|\mathbf{Z}\} \left(1 - \frac{p_k^B}{2}\right) \\ &= \mathbb{E} \left\{ \sum_{k \in \mathcal{A}} W_k \left(1 - \frac{p_k^B}{2}\right) \mathbb{1}_{\{\alpha^A=k\}} + \sum_{k \in \mathcal{B}} \left(E_k - \frac{q_k^B}{2}\right) \mathbb{1}_{\{\alpha^A=k\}} + \sum_{k \in \mathcal{C} \cup \mathcal{A}\mathcal{B}} E_k \left(1 - \frac{p_k^B}{2}\right) \mathbb{1}_{\{\alpha^A=k\}} \middle| \mathbf{Z} \right\}. \end{aligned} \quad (14)$$

The above expectation is maximized when A chooses according to the given policy. ■

Next, we find a potential function for the game. A potential function is a function of the strategies of the players such that the change of the utility of a player when he changes his strategy (while the strategies of other players are held fixed) is equal to the change of the potential function [15].

*Theorem 1:* The function  $H(g^A, g^B)$  given by,

$$H(g^A, g^B) = \sum_{k \in \mathcal{A}} (q_k^A + E_k p_k^B) + \sum_{k \in \mathcal{B}} (q_k^B + E_k p_k^A) + \sum_{k \in \mathcal{C} \cup \mathcal{A} \cup \mathcal{B}} E_k (p_k^A + p_k^B) - \frac{1}{2} \left( \sum_{k \in \mathcal{A}} q_k^A p_k^B + \sum_{k \in \mathcal{B}} p_k^A q_k^B + \sum_{k \in \mathcal{C} \cup \mathcal{A} \cup \mathcal{B}} E_k p_k^A p_k^B \right), \quad (15)$$

is a potential function for the game, where  $p_k^A, p_k^B$  for  $1 \leq k \leq n$ ,  $q_k^A$  for  $k \in \mathcal{A}$  and  $q_k^B$  for  $k \in \mathcal{B}$  are defined in (5) and (6). Moreover, we have that for all  $g^A, g^B \in \mathcal{S}^A \times \mathcal{S}^B$ ,  $H(g^A, g^B) \leq 2 \sum_{k=1}^n E_k$ .

*Proof:* The key to the proof is separating (15) (using (11), (12)) as,

$$H(g^A, g^B) = \mathbb{E}\{R^A(g^A, g^B)\} + \sum_{k \in \mathcal{B}^c} E_k p_k^B + \sum_{k \in \mathcal{B}} q_k^B \quad (16)$$

$$= \mathbb{E}\{R^B(g^A, g^B)\} + \sum_{k \in \mathcal{A}} q_k^A + \sum_{k \in \mathcal{A}^c} p_k^A E_k. \quad (17)$$

Consider updating the strategy of player A while holding the strategy of player B fixed. Notice that since  $\sum_{k \in \mathcal{B}^c} E_k p_k^B + \sum_{k \in \mathcal{B}} q_k^B$  is not affected in this process, from (16), we have that the change in the expected utility of player A is equal to the change of the  $H$  function. Similar holds when player B updates the strategy while holding player A's strategy fixed. Hence, this is indeed a potential function.

To prove the result on the boundedness of  $H(g^A, g^B)$ , notice that from the definition of  $H(g^A, g^B)$ , we have that,

$$\begin{aligned} H(g^A, g^B) &\leq \sum_{k \in \mathcal{A}} (q_k^A + E_k p_k^B) + \sum_{k \in \mathcal{B}} (q_k^B + E_k p_k^A) + \sum_{k \in \mathcal{C} \cup \mathcal{A} \cup \mathcal{B}} E_k (p_k^A + p_k^B) \\ &= \sum_{k \in \mathcal{A}} (\mathbb{E}\{W_k \mathbb{1}_{\alpha^A=k} | \mathbf{Z}\} + E_k \mathbb{E}\{\mathbb{1}_{\alpha^B=k} | \mathbf{Z}\}) + \sum_{k \in \mathcal{B}} (\mathbb{E}\{W_k \mathbb{1}_{\alpha^B=k} | \mathbf{Z}\} + E_k \mathbb{E}\{\mathbb{1}_{\alpha^A=k} | \mathbf{Z}\}) \\ &\quad + \sum_{k \in \mathcal{C} \cup \mathcal{A} \cup \mathcal{B}} E_k (\mathbb{E}\{\mathbb{1}_{\alpha^A=k} | \mathbf{Z}\} + \mathbb{E}\{\mathbb{1}_{\alpha^B=k} | \mathbf{Z}\}) \\ &\leq 2 \sum_{k=1}^n E_k, \end{aligned} \quad (18)$$

where the last inequality follows since  $\mathbb{1}_{\alpha^A=k}, \mathbb{1}_{\alpha^B=k} \leq 1$  and  $W_i$  are independent.  $\blacksquare$

Using Theorem 1 with standard potential game theory (See, for example [36]), we have that the iterative best response algorithm with the best response found in Lemma 1, converges to an  $\epsilon$ -approximate Nash equilibrium in at most  $(2 \sum_{k=1}^n E_k)/\epsilon$  iterations.

## V. WORST-CASE EXPECTED UTILITY

Finding a Nash equilibrium using the above algorithm may not be desirable when the players do not trust each other and place no assumptions on the incentives of the opponent. To mitigate this issue, we consider maximizing the worst-case expected utility of player A. Similar to the case of finding the Nash equilibrium, the analysis is simplified when  $\mathbf{Z}$  is fixed.

Notice that we can simplify (10) to yield,

$$\begin{aligned} & \sum_{k=1}^n \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^A=k\}} \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \\ &= \sum_{k \in \mathcal{A}} q_k^A \mathbb{E}\{\mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} + \sum_{k \in \mathcal{B}} p_k^A \mathbb{E}\{W_k \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} + \sum_{k \in \mathcal{C} \cup \mathcal{AB}} E_k p_k^A \mathbb{E}\{\mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} \\ &= \sum_{k \in \mathcal{A}} \mathbb{E}\{\Omega_k q_k^A \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\} + \sum_{k \in \mathcal{A}^c} \mathbb{E}\{\Omega_k p_k^A \mathbb{1}_{\{\alpha^B=k\}} | \mathbf{Z}\}, \end{aligned} \quad (19)$$

$$\quad (20)$$

where

$$\Omega_k = \begin{cases} 1 & \text{if } k \in \mathcal{A}, \\ W_k & \text{if } k \in \mathcal{B}, \\ E_k & \text{if } k \in \mathcal{C} \cup \mathcal{AB}. \end{cases} \quad (21)$$

Plugging the above in (9), we get that,

$$\mathbb{E}\{R^A(g^A, g^B) | \mathbf{Z}\} = \sum_{k \in \mathcal{A}} q_k^A + \sum_{k \in \mathcal{A}^c} E_k p_k^A - \frac{1}{2} \mathbb{E} \left\{ \sum_{k \in \mathcal{A}} \Omega_k q_k^A \mathbb{1}_{\{\alpha^B=k\}} + \sum_{k \in \mathcal{A}^c} \Omega_k p_k^A \mathbb{1}_{\{\alpha^B=k\}} \middle| \mathbf{Z} \right\}. \quad (22)$$

The difficulty in dealing with  $\mathbb{E}\{R^A(g^A, g^B) | \mathbf{Z}\}$  is that it depends on the strategy  $g^B$  of player B which is not known to player A. Hence, given a strategy  $g^A$  of player A, we first focus on obtaining the worst-case strategy  $\widehat{g^A}$  of player B. Then we focus on finding the strategy  $g^A$  of player A which maximizes  $\mathbb{E}\{R^A(g^A, \widehat{g^A}) | \mathbf{Z}\}$ . This way, we can guarantee a minimum expected utility for player A irrespective of player B's strategy.

*Lemma 2:* For given  $g^A \in \mathcal{S}^A$ , the strategy  $g^B \in \mathcal{S}^B$  that minimizes  $\mathbb{E}\{R^A(g^A, g^B)|\mathbf{Z}\}$  chooses  $\alpha^B = \arg \max_{1 \leq k \leq n} \Lambda_k$ , where,

$$\Lambda_k = \begin{cases} \Omega_k q_k^A & \text{if } k \in \mathcal{A}, \\ \Omega_k p_k^A & \text{if } k \in \mathcal{A}^c, \end{cases} \quad (23)$$

and  $\Omega_k$  are defined in (21).

*Proof:* Notice that the only term of  $\mathbb{E}\{R^A(g^A, g^B)|\mathbf{Z}\}$  in (22) that depends on the strategy of player B is the last expectation. This expectation is maximized when player B chooses  $k$  for which  $\Lambda_k$  is maximized. <sup>1</sup> ■

Hence we have,

$$\mathbb{E}\{R^A(g^A, \widehat{g^A})|\mathbf{Z}\} = \sum_{k \in \mathcal{A}} q_k^A + \sum_{k \in \mathcal{A}^c} E_k p_k^A - \frac{1}{2} \mathbb{E}\{\max\{\Lambda_k; 1 \leq k \leq n\}|\mathbf{Z}\}, \quad (24)$$

where  $\Lambda_k$  are defined in (23). We formulate a strategy for player A using the following optimization problem

$$\begin{aligned} \text{(P1): maximize} \quad & f(\mathbf{q}, \mathbf{p}_{a+1:n}) \\ & g \in \mathcal{S}^A \\ \text{subject to} \quad & \mathbf{q} \in \mathbb{R}^a, \\ & \mathbf{p} \in \mathbb{R}^n, \\ & q_k = \mathbb{E}_{\mathbf{W}, U^A} \{W_k \mathbb{1}_{\{g(U^A, \mathbf{X}, \mathbf{Z})=k\}}|\mathbf{Z}\} \quad \forall 1 \leq k \leq a, \\ & p_l = \mathbb{E}_{\mathbf{W}, U^A} \{\mathbb{1}_{\{g(U^A, \mathbf{X}, \mathbf{Z})=l\}}|\mathbf{Z}\} \quad \forall 1 \leq l \leq n, \end{aligned} \quad (25)$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is defined by,

$$f(\mathbf{x}) = \sum_{k \in \mathcal{A}} x_k + \sum_{k \in \mathcal{A}^c} E_k x_k - \frac{1}{2} \mathbb{E}\{\max\{\Omega_j x_j; 1 \leq j \leq n\}|\mathbf{Z}\}. \quad (26)$$

Although not used immediately, we derive certain properties of  $f$  in the following theorem, which are useful later.

*Theorem 2:* The function  $f$

<sup>1</sup>Ideally, player B may not have information about  $q_j^A$  and  $p_j^A$ . Hence, player B may not be able to utilize this exact strategy. Nevertheless, obtaining a better bound is impossible since we do not have any assumptions or information about player B's strategy. For instance, if player B assumes that player A is using a particular strategy and if player B's assumption turns out to be correct since player B knows the distributions of all  $W_j$  for  $1 \leq j \leq n$ , player B's estimates of  $q_j^A$  and  $p_j^A$  are exact.

- 1) is concave.
- 2) is entry-wise non-decreasing.
- 3) satisfies,

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \frac{3}{2} \sum_{j \in \mathcal{A}} |x_j - y_j| + \frac{3}{2} \sum_{j \in \mathcal{A}^c} E_j |x_j - y_j|, \quad (27)$$

for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ .

*Proof:* See Appendix A. ■

It is important to notice that the values of  $p_k^A$  for  $1 \leq k \leq n$  and  $q_l^A$  for  $l \in \mathcal{A}$  defined in (5) completely determine the optimal value of (P1). Hence notice that if we have a mechanism to find the set  $\mathcal{G}^A \subset \mathbb{R}^{n+a}$  defined by,

$$\mathcal{G}^A = \left\{ (\mathbf{q}, \mathbf{p}) \left| \begin{array}{l} g \in \mathcal{S}^A, \mathbf{p} \in \mathbb{R}^n, p_k = \mathbb{E}\{\mathbb{1}_{\{g(U^A, \mathbf{X}, \mathbf{Z})=k\}} | \mathbf{Z}\} \text{ for } 1 \leq k \leq n, \\ \mathbf{q} \in \mathbb{R}^a, q_l = \mathbb{E}\{W_l \mathbb{1}_{\{g(U^A, \mathbf{X}, \mathbf{Z})=l\}} | \mathbf{Z}\} \text{ for } 1 \leq l \leq a \end{array} \right. \right\}, \quad (28)$$

we can solve the optimization problem,

$$\begin{aligned} \text{(P1.1): maximize} \quad & f(\mathbf{q}, \mathbf{p}_{a+1:n}) \\ & (\mathbf{q}, \mathbf{p}) \\ \text{subject to} \quad & (\mathbf{q}, \mathbf{p}) \in \mathcal{G}^A, \end{aligned} \quad (29)$$

to find the optimal  $(\mathbf{p}^*, \mathbf{q}^*)$ , after which we find the optimal strategy  $g^*$  for player A that satisfies

$$\begin{aligned} p_k^* &= \mathbb{E}\{\mathbb{1}_{\{g^*(U^A, \mathbf{X}, \mathbf{Z})=k\}} | \mathbf{Z}\} \text{ for } 1 \leq k \leq n, \\ q_l^* &= \mathbb{E}\{W_l \mathbb{1}_{\{g^*(U^A, \mathbf{X}, \mathbf{Z})=l\}} | \mathbf{Z}\} \text{ for } 1 \leq l \leq a. \end{aligned} \quad (30)$$

It turns out that  $\mathcal{G}^A$  is a convex set, as established by the following lemma. Hence combining with Theorem 2, we have that (P1.1) is a convex optimization problem.

*Theorem 3:*  $\mathcal{G}^A$  is a convex set. Further in the special case  $a = 0$ ,  $\mathcal{G}^A$  is the  $n$ -dimensional probability simplex  $\mathcal{I} = \{\mathbf{p} \in \mathbb{R}^n : \sum_{i=1}^n p_i = 1, p_i \geq 0 \forall i\}$ .

*Proof:* See Appendix B. ■

It turns out that finding  $\mathcal{G}^A$  for the general scenario is difficult, although the case  $a = 0$  is handled by Theorem 3. In fact, when  $a = b = d = 0$ , an explicit solution can be obtained to (P1.1), which we describe in section V-A. In section V-B, we describe the solution to the general case. In Appendix H, we provide simpler alternative solutions to the special cases  $a = 0$  (with no restriction on  $b$ ) and  $a = 1$  (with the additional assumption that  $W_1$  has a continuous CDF).

A. *Explicit solution for  $a = b = d = 0$*

When neither player knows any of the reward realizations, we have  $a = b = d = 0$ , and the problem reduces to the following.

$$\begin{aligned} \text{(P2): maximize} \quad & \sum_{k=1}^n p_k E_k - \frac{1}{2} \max\{p_k E_k; 1 \leq k \leq n\} \\ \text{subject to} \quad & \mathbf{p} \in \mathcal{I}, \end{aligned} \tag{31}$$

where,

$$\mathcal{I} = \{\mathbf{p} \in \mathbb{R}^n : \sum_{i=1}^n p_i = 1, p_i \geq 0 \forall i\} \tag{32}$$

is the  $n$ -dimensional probability simplex. For this section, we assume without loss of generality that  $E_k > 0$  for all  $k$ . If at least one of the  $E_k$ 's were zero, we could transform (P2) into a lower dimensional problem with non-zero  $E_k$ 's. The following lemma constructs an explicit solution for  $a = b = d = 0$ .<sup>2</sup>

*Lemma 3:* Assume without loss of generality that  $E_k \geq E_{k+1}$  for  $1 \leq k \leq n - 1$ . Also, let,

$$r = \arg \max_{1 \leq k \leq n} \frac{k - \frac{1}{2}}{\sum_{j=1}^k \frac{1}{E_j}}, \tag{33}$$

where the lowest index is chosen in the case of ties. The optimal solution for (P2) is given by  $\mathbf{p}^*$  where,

$$p_k^* = \begin{cases} \frac{1}{E_k \left( \sum_{j=1}^r \frac{1}{E_j} \right)} & \text{if } k \leq r, \\ 0 & \text{otherwise.} \end{cases} \tag{34}$$

*Proof:* See Appendix C. ■

It should be noted that this solution is not unique. For instance, consider the case when  $n = 2$ ,  $E_1 = 2$ , and  $E_2 = 1$ . In this case, the lemma finds the solution  $(p_1, p_2) = (1, 0)$ , but it should be noted that  $(p_1, p_2) = (1/3, 2/3)$  is also a solution. It is also interesting that the solution assigns positive probabilities to the  $r$  resources with the highest average reward, although within these  $r$  resources, higher probabilities are assigned to the resources with lower rewards.

<sup>2</sup>The same problem structure arises in the case with symmetric information between the players (case  $a = b = 0$  with  $d$  arbitrary). Hence, we can use the solution obtained in this section for the above case as well.

It should also be noted that the worst-case strategy can be arbitrarily worse than the Nash equilibrium strategy. For instance, consider the simple scenario with two resources such that  $E_1 = E_2$ , where none of the players observe any of the reward realizations. In this case, a Nash equilibrium would be player A always choosing resource 1 and player B always choosing resource 2. Another Nash equilibrium would be player B always choosing resource 1 and player A always choosing resource 2. In either case, player A's expected utility is  $E_1$ . However, notice that, from Lemma 3, the maximum worst-case expected utility of player A is  $3E_1E_2/(2E_1 + 2E_2) = 3E_1/4$ . Hence,  $E_1$  can be scaled to obtain arbitrarily large deviation between the worst-case and the Nash equilibrium solutions.

### B. Solving the general case

In this section, we focus on solving the most general version of (P1) (with no restrictions on the sets  $\mathcal{A}, \mathcal{B}, \mathcal{AB}, \mathcal{C}$ ). In particular, we focus on finding a mixed strategy to optimize the worst-case expected utility for player A. It turns out that our optimal solution chooses from a mixture of pure strategies parameterized by  $\mathbf{Q} \in \mathbb{R}^n$ , of the following form

$$g_{\mathbf{Q}}^A(\mathbf{X}) = \arg \max_{1 \leq j \leq n} \{ \{Q_j W_j; j \in \mathcal{A}\} \cup \{Q_j; j \in \mathcal{A}^c\} \}. \quad (35)$$

We name this special class of pure strategies as *threshold strategies*. We develop a novel algorithm to solve this problem. Our algorithm leverages techniques from drift-plus-penalty theory [14] and online convex optimization [12], [13]. It should be noted that our algorithm runs offline and is used to construct an appropriate strategy for player A that approximately solves (P1) conditioned on the observed realization of  $\mathbf{Z}$ . We show that we can get arbitrarily close to the optimal value of (P1) by using a finite equiprobable mixture of pure strategies of the above form. It should be noted that the algorithm developed in this section can be used to solve the general unconstrained problem of finding the randomized decision  $\alpha \in \{1, 2, \dots, n\}$  which maximizes  $\mathbb{E}\{h(\mathbf{x}; \Theta)\}$ , where  $\mathbf{x} \in \mathbb{R}^n$  with  $x_k = \mathbb{E}\{\Gamma_k \mathbb{1}_{\{\alpha=k\}}\}$ ,  $\Theta \in \mathbb{R}^m$  and  $\Gamma \in \mathbb{R}^n$  are non-negative random vectors with finite second moments, and  $h$  is a concave function such that  $\tilde{h}(\mathbf{x}) = \mathbb{E}\{h(\mathbf{x}; \Theta)\}$  is Lipschitz continuous, entry-wise non-decreasing and has bounded subgradients.

We first provide an algorithm that generates a mixture of  $T$  pure strategies, after which we establish the closeness to the optimality of the mixture. We generate the mixture of  $T$  pure strategies  $\{g_{\mathbf{Q}(t)}^A\}_{t=1}^T$  by iteratively updating vector  $\mathbf{Q}$  for  $T$  iterations, where  $\mathbf{Q}(t)$ , and

$g_{\mathbf{Q}(t)}(\mathbf{X})$  denote the state of  $\mathbf{Q}$  and the pure strategy generated in the  $t$ -th iteration, respectively. In addition to  $\mathbf{Q}(t)$ , we require another state vector  $\boldsymbol{\gamma}(t) \in \mathbb{R}^n$ , which we also update in each iteration, and a parameter  $V$  which decides the convergence properties of the algorithm. We provide the specific details on setting  $V$  later in our analysis. We begin with  $\mathbf{Q}(1) = \boldsymbol{\gamma}(0) = 0$ . In the  $t$ -th iteration ( $t \geq 1$ ), we independently sample  $\mathbf{X}(t)$  and  $\boldsymbol{\Omega}(t)$  from the distributions of  $\mathbf{X}$  and  $\boldsymbol{\Omega}$ , respectively, where  $\boldsymbol{\Omega}$  is defined in (21), while keeping  $\mathbf{Z}$  fixed to its observed value. Then we update  $\boldsymbol{\gamma}(t)$  and  $\mathbf{Q}(t+1)$  as follows. First, we solve,

$$\text{(P3): minimize}_{\boldsymbol{\gamma}(t)} \quad -V f'_t(\boldsymbol{\gamma}(t-1))^\top \boldsymbol{\gamma}(t) + \alpha \|\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\|_2^2 + \sum_{j=1}^n Q_j(t) \gamma_j(t) \quad (36a)$$

$$\text{subject to} \quad \boldsymbol{\gamma}(t) \in \mathcal{K}, \quad (36b)$$

to find  $\boldsymbol{\gamma}(t)$ , where

$$f_t(\mathbf{x}) = \sum_{k \in \mathcal{A}} x_k + \sum_{k \in \mathcal{A}^c} x_k E_k - \frac{1}{2} \max\{x_k \Omega_k(t); 1 \leq k \leq n\}, \quad (37)$$

$\alpha > 0$  and  $\mathcal{K} = \left(\prod_{j \in \mathcal{A}} [0, E_j]\right) \times [0, 1]^{n-a}$ . Notice that  $f'_t(\mathbf{x})$  is given by,

$$f'_{t,j}(\mathbf{x}) = \begin{cases} 1 - \frac{1}{2} \mathbb{1}_{\{\arg \max_{1 \leq k \leq n} \{x_k \Omega_k(t)\} = j\}} & \text{if } j \in \mathcal{A}, \\ E_j - \frac{1}{2} \mathbb{1}_{\{\arg \max_{1 \leq k \leq n} \{x_k \Omega_k(t)\} = j\}} \Omega_j(t) & \text{if } j \in \mathcal{A}^c, \end{cases} \quad (38)$$

where  $\arg \max$  returns the lowest index in the case of ties. Notice that  $f_t$  is a concave function, which can be established by repeating the same argument used to establish the concavity of  $f$  in Theorem 2. Then we choose the action for the  $t$ -th iteration  $\alpha^A(t) = g_{\mathbf{Q}(t)}^A(\mathbf{X}(t))$  (See (35)).

Then to update  $\mathbf{Q}(t+1)$ , we use,

$$\begin{aligned} Q_j(t+1) &= \max\{Q_j(t) + \gamma_j(t) - X_j(t) \mathbb{1}_{\{\alpha^A(t)=j\}}, 0\}, \forall j \in \mathcal{A}, \\ Q_j(t+1) &= \max\{Q_j(t) + \gamma_j(t) - \mathbb{1}_{\{\alpha^A(t)=j\}}, 0\}, \forall j \in \mathcal{A}^c. \end{aligned} \quad (39)$$

The algorithm is summarized as Algorithm 1 for clarity.

After creating the mixture  $\{g_{\mathbf{Q}(t)}^A\}_{t=1}^T$  of pure strategies, we choose one of them randomly with probability  $1/T$  to take the decision. In the following two sections, we focus on solving (P3) and evaluating the performance of the Algorithm 1.



---

**Algorithm 1:** Algorithm to generate the optimal mixture of  $T$  pure strategies

---

```

1 Initialize  $\mathbf{Q}(1) = \boldsymbol{\gamma}(0) = 0$ 
2 for each iteration  $t \in [1 : T]$  do
3   | Sample  $\mathbf{X}(t)$ , and  $\boldsymbol{\Omega}(t)$ 
4   | Choose  $\boldsymbol{\gamma}(t)$  by solving (P3)
5   | Choose the action  $\alpha^A(t) = g_{\mathbf{Q}(t)}^A(\mathbf{X}(t))$ 
6   | Obtain  $\mathbf{Q}(t+1)$  using (39)
7 end

```

---

1) *Solving (P3):* Notice that the objective of (P3) can be written as

$$\begin{aligned}
& -V f'_t(\boldsymbol{\gamma}(t-1))^\top \boldsymbol{\gamma}(t) + \alpha \|\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\|_2^2 + \sum_{j=1}^n Q_j(t) \gamma_j(t) \\
& = \sum_{j=1}^n \left\{ -V f'_{t,j}(\boldsymbol{\gamma}(t-1)) \gamma_j(t) + \alpha (\gamma_j(t) - \gamma_j(t-1))^2 + Q_j(t) \gamma_j(t) \right\}. \tag{40}
\end{aligned}$$

Hence (P3) seeks to minimize a separable convex function over the box constraint  $\boldsymbol{\gamma}(t) \in \mathcal{K}$ . The solution vector  $\boldsymbol{\gamma}(t)$  is found by separately minimizing each component  $\gamma_j(t)$  over  $[0, u_j]$ , where

$$u_j = \begin{cases} E_j & \text{if } j \in \mathcal{A}, \\ 1 & \text{if } j \in \mathcal{A}^c. \end{cases} \tag{41}$$

The resulting solution is,

$$\gamma_j(t) = \Pi_{[0, u_j]} \left( \gamma_j(t-1) - \frac{-V f'_{t,j}(\boldsymbol{\gamma}(t-1)) + Q_j(t)}{2\alpha} \right), \tag{42}$$

where  $\Pi_{[0, u_j]}$  denotes the projection onto  $[0, u_j]$ . Notice that the above solution is obtained by projecting the global minimizer of the function to be minimized onto  $[0, u_j]$ .

2) *How good is the mixed strategy generated by Algorithm 1:* Without loss of generality, we assume that  $E_k > 0$  for all  $1 \leq k \leq n$ . The following theorem establishes the closeness of the expected utility generated by Algorithm 1 to the optimal value  $f^{\text{opt}}$  of (P1).

*Theorem 4:* Assume  $\alpha$  is set such that  $\alpha \geq V^2$ , and we use the mixed strategy  $g^A$  generated by Algorithm 1 to make the decision. Then,

$$\begin{aligned} \mathbb{E}\{R^A(g^A, \widehat{g^A})|\mathbf{Z}\} \geq & f^{\text{opt}} - \frac{D_1}{V} - \frac{VD_2}{16\alpha} - \frac{\alpha D_3}{VT} - \frac{3}{2T} \sum_{k \in \mathcal{A}} \left\{ \sqrt{\alpha} + E_k \left( 2\sqrt{2\alpha} + 1 \right) \right\} \\ & - \frac{3}{2T} \sum_{k \in \mathcal{A}^c} \left\{ E_k^2 \sqrt{\alpha} + E_k \left( 2\sqrt{2\alpha} + 1 \right) \right\}, \end{aligned} \quad (43)$$

where,

$$\begin{aligned} D_1 &= n - a + \frac{1}{2} \sum_{j \in \mathcal{A}} (E_j^2 + \mathbb{E}\{W_k^2\}), \\ D_2 &= 4a + \mathbb{E}\{\|\boldsymbol{\Omega}\|_2^2|\mathbf{Z}\} + \sum_{j \in \mathcal{A}^c} 4E_j^2, \\ D_3 &= n - a + \sum_{j \in \mathcal{A}} E_j^2, \end{aligned} \quad (44)$$

$\boldsymbol{\Omega}$  is defined in (21), and  $f^{\text{opt}}$  is the optimal value of (P1). Hence, by fixing  $\varepsilon > 0$ , and using  $V = 1/\varepsilon$ ,  $\alpha = 1/\varepsilon^2$ , and  $T \geq 1/\varepsilon^2$ , the average error is  $\mathcal{O}(\varepsilon)$ .

*Proof:* The key to the proof is noticing that  $\mathbf{Q}(t)$  can be treated as  $n$  queues. Before proceeding with the proof, we define some quantities. Define the history up to time  $t$  by  $\mathcal{H}(t) = \{\mathbf{X}(\tau); 1 \leq \tau < t\} \cup \{\boldsymbol{\Omega}(\tau); 1 \leq \tau \leq t\}$ . Notice that we include  $\boldsymbol{\Omega}(t)$  in  $\mathcal{H}(t)$  since this will allow us to treat  $\boldsymbol{\gamma}(t)$ , and  $\mathbf{Q}(t)$  as deterministic functions of  $\mathcal{H}(t)$  and  $\mathbf{Z}$ . Let us define the Lyapunov function  $L(t) = \frac{1}{2} \|\mathbf{Q}(t)\|^2 = \frac{1}{2} \sum_{j=1}^n Q_j(t)^2$ , and the drift  $\Delta(t) = \mathbb{E}\{L(t+1) - L(t)|\mathcal{H}(t), \mathbf{Z}\}$ . Now notice that,

$$\mathbb{E}\{R^A(g^A, \widehat{g^A})|\mathbf{Z}\} = f \left( \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{\mathbf{x}(t)|\mathbf{Z}\} \right), \quad (45)$$

where,

$$x_k(t) = \begin{cases} X_k(t) \mathbb{1}_{\{g_{\mathbf{Q}(t)}^A(\mathbf{x}(t))=k\}} & \text{if } k \in \mathcal{A}, \\ \mathbb{1}_{\{g_{\mathbf{Q}(t)}^A(\mathbf{x}(t))=k\}} & \text{if } k \in \mathcal{A}^c. \end{cases} \quad (46)$$

We begin with the following two lemmas, which will be useful in the proof.

*Lemma 4:* The drift is bounded above as,

$$\Delta(t) \leq D_1 + \sum_{j=1}^n Q_j(t) (\gamma_j(t) - \mathbb{E}\{x_j(t)|\mathcal{H}(t), \mathbf{Z}\}), \quad (47)$$

where  $D_1$  is defined in (44).

*Proof:* See Appendix D. ■

The following is a well-known result regarding the minimization of strongly convex functions (See, for example [37]).

*Lemma 5:* Let  $\mathcal{G} \subset \mathbb{R}^n$  be a convex set with a non-empty interior  $\mathcal{G}^0$ . Let  $\mathcal{C} \subset \mathcal{G}$  such that  $\mathcal{C}$  intersects  $\mathcal{G}^0$ . Define  $\mathcal{C}^0 = \mathcal{C} \cap \mathcal{G}^0$ . Let  $\omega : \mathcal{G} \rightarrow \mathbb{R}$  be a function that is continuously differentiable in  $\mathcal{G}^0$ . The Bregman divergence  $D_\omega : \mathcal{C} \times \mathcal{C}^0 \rightarrow \mathbb{R}$  generated by  $\omega$  is denoted by,

$$D_\omega(\mathbf{x} \parallel \mathbf{y}) = \omega(\mathbf{x}) - \omega(\mathbf{y}) - \nabla \omega(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}). \quad (48)$$

Let  $g : \mathcal{G} \rightarrow \mathbb{R}$  be a convex function. Fix  $\alpha > 0$ , and  $\mathbf{y} \in \mathcal{C}^0$ . Let,

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathcal{C}} g(\mathbf{x}) + \alpha D_\omega(\mathbf{x} \parallel \mathbf{y}). \quad (49)$$

Additionally, assume that  $\mathbf{x}^* \in \mathcal{C}^0$ . Then,

$$g(\mathbf{x}^*) + \alpha D_\omega(\mathbf{x}^* \parallel \mathbf{y}) \leq g(\mathbf{z}) + \alpha D_\omega(\mathbf{z} \parallel \mathbf{y}) - \alpha D_\omega(\mathbf{z} \parallel \mathbf{x}^*), \quad (50)$$

for all  $\mathbf{z} \in \mathcal{C}$ . We have the following special cases.

- Using  $\mathcal{G} = \mathbb{R}^n$ , we have that  $\mathcal{C} = \mathcal{C}^0$ . In this case, if  $\omega(\mathbf{x}) = \|\mathbf{x}\|_2^2$ , we have  $D_\omega(\mathbf{x} \parallel \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ .
- Using  $\mathcal{G} = [0, \infty)^n$  and  $\mathcal{C} = \mathcal{I}$ , we have  $\mathcal{C}^0 = \mathcal{I}^0$ . In this case, if  $\omega(\mathbf{x}) = \sum_{j=1}^n x_j \ln(x_j)$ , we have that,  $D_\omega(\mathbf{x} \parallel \mathbf{y}) = D(\mathbf{x} \parallel \mathbf{y})$ , where  $D(\mathbf{x} \parallel \mathbf{y})$  is the Kullback-Leibler divergence.

Now we move on to the main proof. Notice that the objective of (P3) can be written as

$$g'_t(\boldsymbol{\gamma}(t-1))^\top \boldsymbol{\gamma}(t) + \alpha \|\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\|_2^2, \quad (51)$$

where,

$$g_t(\mathbf{x}) = -V f_t(\mathbf{x}) + \sum_{j=1}^n Q_j(t) x_j. \quad (52)$$

Let  $g^{A,*}$  be the strategy that is optimal for (P1). Let us define  $\mathbf{x}^*(t) \in \mathbb{R}^n$ , where

$$x_k^*(t) = \begin{cases} X_k(t) \mathbb{1}_{\{g^{A,*}(U^A(t), \mathbf{X}(t), \mathbf{Z})=k\}} & \text{if } k \in \mathcal{A}, \\ \mathbb{1}_{\{g^{A,*}(U^A(t), \mathbf{X}(t), \mathbf{Z})=k\}} & \text{if } k \in \mathcal{A}^c, \end{cases} \quad (53a)$$

where  $U^A(t)$  for  $1 \leq t \leq T$  is a collection of i.i.d uniform-[0, 1) random variables. Notice that  $\mathbf{y}^* = \mathbb{E}\{\mathbf{x}^*(t) | \mathbf{Z}\}$  is independent of  $t$  and belongs to  $\mathcal{K}$ . Hence  $\mathbf{y}^*$  is feasible for (P3). Notice that,

$$-V f'_t(\boldsymbol{\gamma}(t-1))^\top \boldsymbol{\gamma}(t) + \sum_{j=1}^n Q_j(t) \gamma_j(t) + \alpha \|\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\|_2^2$$

$$\begin{aligned}
&= g'_t(\boldsymbol{\gamma}(t-1))^\top \boldsymbol{\gamma}(t) + \alpha \|\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\|_2^2 \\
&\leq_{(a)} g'_t(\boldsymbol{\gamma}(t-1))^\top \mathbf{y}^* + \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t-1)\|_2^2 - \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t)\|_2^2 \\
&= -V f'_t(\boldsymbol{\gamma}(t-1))^\top \mathbf{y}^* + \sum_{j=1}^n Q_j(t) y_j^* + \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t-1)\|_2^2 - \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t)\|_2^2, \quad (54)
\end{aligned}$$

where (a) follows from Lemma 5 for the convex function  $h$  given by  $h(\mathbf{x}) = g'_t(\boldsymbol{\gamma}(t-1))^\top \mathbf{x}$  with  $\mathcal{G} = \mathbb{R}^n$ , and  $\mathcal{C} = \mathcal{K}$ , since  $\boldsymbol{\gamma}(t)$  is the solution to (P3) and  $\mathbf{y}^*$  is feasible for (P3). Also, step 5 in each iteration of Algorithm 1 of finding the action can be represented as the maximization of

$$\sum_{j \in \mathcal{A}} Q_j(t) \mathbb{E}\{X_j(t) \mathbb{1}_{\{\alpha^A=j\}} | \mathcal{H}(t), \mathbf{Z}\} + \sum_{j \in \mathcal{A}^c} Q_j(t) \mathbb{E}\{\mathbb{1}_{\{\alpha^A=j\}} | \mathcal{H}(t), \mathbf{Z}\} \quad (55)$$

over all possible actions  $\alpha^A \in \{1, 2, \dots, n\}$  at time-slot  $t$ . Hence comparing the scenario where  $g_{Q(t)}^A$  is used in the  $t$ -th iteration with the scenario where  $g^{A,*}$  is used with randomization variable  $U^A(t)$  in the  $t$ -th iteration, we have the inequality,

$$-\sum_{j=1}^n Q_j(t) \mathbb{E}\{x_j(t) | \mathcal{H}(t), \mathbf{Z}\} \leq -\sum_{j=1}^n Q_j(t) \mathbb{E}\{x_j^*(t) | \mathcal{H}(t), \mathbf{Z}\} = -\sum_{j=1}^n Q_j(t) y_j^*, \quad (56)$$

where the last equality follows since  $\mathbf{x}^*(t)$  is independent of  $\mathcal{H}(t)$  conditioned on  $\mathbf{Z}$ . Summing (54) and (56),

$$\begin{aligned}
&-V f'_t(\boldsymbol{\gamma}(t-1))^\top \boldsymbol{\gamma}(t) + \alpha \|\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\|_2^2 + \sum_{j=1}^n Q_j(t) (\gamma_j(t) - \mathbb{E}\{x_j(t) | \mathcal{H}(t), \mathbf{Z}\}) \quad (57) \\
&\leq -V f'_t(\boldsymbol{\gamma}(t-1))^\top \mathbf{y}^* + \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t-1)\|_2^2 - \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t)\|_2^2.
\end{aligned}$$

Adding  $D_1 + V f'_t(\boldsymbol{\gamma}(t-1))^\top \boldsymbol{\gamma}(t-1)$  to both sides and using Lemma 4 yields,

$$\begin{aligned}
&\Delta(t) - V f'_t(\boldsymbol{\gamma}(t-1))^\top \{\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\} + \alpha \|\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\|_2^2 \\
&\leq D_1 - V f'_t(\boldsymbol{\gamma}(t-1))^\top (\mathbf{y}^* - \boldsymbol{\gamma}(t-1)) + \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t-1)\|_2^2 - \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t)\|_2^2 \\
&\leq D_1 - V \{f_t(\mathbf{y}^*) - f_t(\boldsymbol{\gamma}(t-1))\} + \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t-1)\|_2^2 - \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t)\|_2^2, \quad (58)
\end{aligned}$$

where the last inequality follows from the sub-gradient inequality for the concave function  $f_t$ . Now we introduce the following lemma.

*Lemma 6:* We have,

$$-V f'_t(\boldsymbol{\gamma}(t-1))^\top \{\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\} + \alpha \|\boldsymbol{\gamma}(t) - \boldsymbol{\gamma}(t-1)\|_2^2 \geq -\frac{V^2}{4\alpha} \left( a + \sum_{j \in \mathcal{A}^c} E_j^2 \right)$$

$$- \frac{V^2}{16\alpha} \|\boldsymbol{\Omega}(t)\|_2^2, \quad (59)$$

*Proof:* See Appendix E. ■

Substituting the bound from Lemma 6 in (58) we have that,

$$\begin{aligned} \Delta(t) &- \frac{V^2}{4\alpha} \left( a + \sum_{j \in \mathcal{A}^c} E_j^2 \right) - \frac{V^2}{16\alpha} \|\boldsymbol{\Omega}(t)\|_2^2 \\ &\leq D_1 - V \{f_t(\mathbf{y}^*) - f_t(\boldsymbol{\gamma}(t-1))\} + \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t-1)\|_2^2 - \alpha \|\mathbf{y}^* - \boldsymbol{\gamma}(t)\|_2^2, \end{aligned} \quad (60)$$

The above holds for each  $t \in \{1, 2, \dots, T\}$ . Hence we first take the expectation conditioned on  $\mathbf{Z}$  of both sides of the above expression, after which we sum from 1 to  $T$ , which results in,

$$\begin{aligned} &\mathbb{E}\{L(T+1)|\mathbf{Z}\} - \mathbb{E}\{L(1)|\mathbf{Z}\} - \frac{TV^2}{4\alpha} \left( a + \sum_{j \in \mathcal{A}^c} E_j^2 \right) - \frac{TV^2}{16\alpha} \mathbb{E}\{\|\boldsymbol{\Omega}\|_2^2|\mathbf{Z}\} \\ &\leq D_1 T - V \sum_{t=1}^T \mathbb{E}\{f_t(\mathbf{y}^*)|\mathbf{Z}\} + V \sum_{t=1}^T \mathbb{E}\{f_t(\boldsymbol{\gamma}(t-1))|\mathbf{Z}\} + \alpha \mathbb{E}\{\|\mathbf{y}^* - \boldsymbol{\gamma}(0)\|_2^2|\mathbf{Z}\} \\ &\quad - \alpha \mathbb{E}\{\|\mathbf{y}^* - \boldsymbol{\gamma}(T)\|_2^2|\mathbf{Z}\}. \end{aligned} \quad (61)$$

Notice that,

$$\mathbb{E}\{f_t(\mathbf{y}^*)|\mathbf{Z}\} = f(\mathbf{y}^*) = f^{\text{opt}}, \quad (62)$$

where functions  $f$  and  $f_t$  are defined in (26), and (37), respectively. Also, we have that,

$$\mathbb{E}\{f_t(\boldsymbol{\gamma}(t-1))|\mathbf{Z}\} = \mathbb{E}\{\mathbb{E}_{\boldsymbol{\Omega}(t)}\{f_t(\boldsymbol{\gamma}(t-1))|\mathcal{H}(t-1), \mathbf{Z}\}|\mathbf{Z}\} =_{(a)} \mathbb{E}\{f(\boldsymbol{\gamma}(t-1))|\mathbf{Z}\}, \quad (63)$$

where (a) follows from the definition of  $f_t$  in (37), since  $\boldsymbol{\gamma}(t-1)$  is a function of  $\mathcal{H}(t-1)$  and  $\boldsymbol{\Omega}(t)$  is independent of  $\mathcal{H}(t-1)$  conditioned on  $\mathbf{Z}$ . Substituting (62) and (63) in (61), we have that,

$$\begin{aligned} &\mathbb{E}\{L(T+1)|\mathbf{Z}\} - \mathbb{E}\{L(1)|\mathbf{Z}\} - \frac{TV^2}{4\alpha} \left( a + \sum_{j \in \mathcal{A}^c} E_j^2 \right) - \frac{TV^2}{16\alpha} \mathbb{E}\{\|\boldsymbol{\Omega}\|_2^2|\mathbf{Z}\} \\ &\leq D_1 T - VT f^{\text{opt}} + V \sum_{t=1}^T \mathbb{E}\{f(\boldsymbol{\gamma}(t-1))|\mathbf{Z}\} + \alpha \mathbb{E}\{\|\mathbf{y}^* - \boldsymbol{\gamma}(0)\|_2^2|\mathbf{Z}\} - \alpha \mathbb{E}\{\|\mathbf{y}^* - \boldsymbol{\gamma}(T)\|_2^2|\mathbf{Z}\} \\ &\leq_{(a)} D_1 T - VT f^{\text{opt}} + V \sum_{t=1}^T \mathbb{E}\{f(\boldsymbol{\gamma}(t-1))|\mathbf{Z}\} + \alpha \left( n - a + \sum_{k \in \mathcal{A}} E_k^2 \right) \\ &\leq D_1 T - VT f^{\text{opt}} + VT f \left( \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{\boldsymbol{\gamma}(t-1)|\mathbf{Z}\} \right) + \alpha D_3, \end{aligned} \quad (64)$$

where (a) follows since  $\mathbf{y}^*, \gamma(T), \gamma(0) \in \mathcal{K}$ , and the last inequality follows from Jensen's inequality on the concave function  $f$ . (See the definition of  $D_3$  in (44)). Since  $\mathbf{Q}(1) = 0$  and  $\mathbb{E}\{L(T+1)|\mathbf{Z}\} \geq 0$ , after some rearrangements above translates to,

$$f^{\text{opt}} - \frac{D_1}{V} - \frac{VD_2}{16\alpha} - \frac{\alpha D_3}{VT} \leq f\left(\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\{\gamma(t)|\mathbf{Z}\}\right), \quad (65)$$

where  $D_2$  is defined in (44). Now we prove that prove the following lemma.

*Lemma 7:* We have,

$$\begin{aligned} f\left(\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\{\gamma(t)|\mathbf{Z}\}\right) &\leq f\left(\frac{1}{T} \sum_{t=1}^T \mathbb{E}\{\mathbf{x}(t)|\mathbf{Z}\}\right) + \frac{3}{2T} \sum_{k \in \mathcal{A}} \left\{ \sqrt{\alpha} + E_k \left(2\sqrt{2\alpha} + 1\right) \right\} \\ &\quad + \frac{3}{2T} \sum_{k \in \mathcal{A}^c} \left\{ E_k^2 \sqrt{\alpha} + E_k \left(2\sqrt{2\alpha} + 1\right) \right\}. \end{aligned} \quad (66)$$

*Proof:* We first introduce the following two lemmas.

*Lemma 8:* The queues  $Q_j(t)$  for  $1 \leq j \leq n$  updated according to Algorithm 1 satisfy,

$$\max \left\{ \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\gamma_j(t) - x_j(t)|\mathbf{Z}\}, \mathbf{0} \right\} \leq \frac{\mathbb{E}\{Q_j(T) | \mathbf{Z}\}}{T}. \quad (67)$$

*Proof:* See Appendix F. ■

The following lemma is vital in constructing the  $\mathcal{O}(\sqrt{\alpha})$  bound on the queue sizes, which leads to the  $\mathcal{O}(1/\varepsilon^2)$  solution. It should be noted that an easier bound can be obtained on the queue sizes, which leads to a  $\mathcal{O}(1/\varepsilon^3)$  solution.

*Lemma 9:* Given that  $\alpha \geq V^2$ ,  $\mathbf{Q}(t)$  satisfy the bound

$$Q_j(t) \leq \begin{cases} (1 + 2\sqrt{2}E_j)\sqrt{\alpha} + E_j & \text{if } j \in \mathcal{A} \\ (E_j + 2\sqrt{2})\sqrt{\alpha} + 1 & \text{if } j \in \mathcal{A}^c, \end{cases} \quad (68)$$

for each  $t \in [1 : T]$ .

*Proof:* See Appendix G. ■

Now we move on to the main proof. Notice that,

$$\begin{aligned} f\left(\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\{\gamma(t)|\mathbf{Z}\}\right) &= f\left(\frac{\gamma(0)}{T} + \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\gamma(t)|\mathbf{Z}\}\right) \\ &\leq f\left(\frac{1}{T} \sum_{t=1}^T \mathbb{E}\{\mathbf{x}(t)|\mathbf{Z}\} + \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\gamma(t)|\mathbf{Z}\} - \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\mathbf{x}(t)|\mathbf{Z}\} - \frac{\mathbb{E}\{\mathbf{x}(T)|\mathbf{Z}\}}{T}\right) \end{aligned}$$

$$\begin{aligned}
&\leq_{(a)} f \left( \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{\mathbf{x}(t)|\mathbf{Z}\} + \max \left\{ \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\gamma(t)|\mathbf{Z}\} - \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\mathbf{x}(t)|\mathbf{Z}\}, \mathbf{0} \right\} \right) \quad (69) \\
&\leq_{(b)} f \left( \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{\mathbf{x}(t)|\mathbf{Z}\} \right) + \frac{3}{2} \sum_{k \in \mathcal{A}} \left( \max \left\{ \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\gamma_k(t) - x_k(t)|\mathbf{Z}\}, \mathbf{0} \right\} \right) \\
&\quad + \frac{3}{2} \sum_{k \in \mathcal{A}^c} E_k \left( \max \left\{ \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\gamma_k(t) - x_k(t)|\mathbf{Z}\}, \mathbf{0} \right\} \right),
\end{aligned}$$

where (a) follows from the entry-wise non-decreasing property of  $f$  (Theorem 2-2) and (b) follows from Theorem 2-3. Combining (69), and Lemma 8 with the bound on  $\mathbf{Q}(T)$  given by Lemma 9, we are done with the proof of the lemma. ■

Combining Lemma 7 with (65), we are done with the proof of the theorem. ■

## VI. SIMULATIONS

For the simulations, we use  $W_j$  as exponential random variables. Notice that since we are conditioning on  $\mathbf{Z}$  to solve the problem, the objective of (P1) defined in (26) has the same structure for the two scenarios  $(a, b, c, d)$  and  $(a, b, c + d, 0)$ . Hence, we use  $d = 0$  for all the simulations. Notice that the sets  $\mathcal{A}$  and  $\mathcal{B}$  denote the private information of players A and B, respectively. We consider the three scenarios given below.

- 1)  $a = 0, b = 0, c = 3, d = 0$ : Both players do not have private information.
- 2)  $a = 0, b = 1, c = 2, d = 0$ : Only player B has private information.
- 3)  $a = 1, b = 1, c = 1, d = 0$ : Both players have private information.

Figures 1–3 show pictorial representations of these cases.

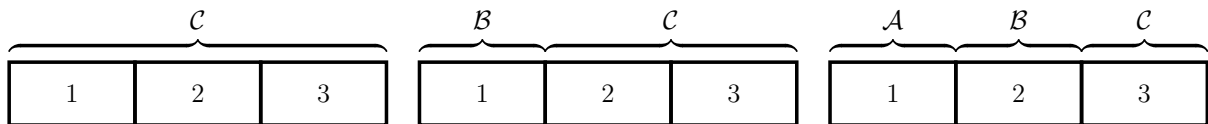


Fig. 1. Case  $a, b, c, d = 0, 0, 3, 0$

Fig. 2. Case  $a, b, c, d = 0, 1, 2, 0$

Fig. 3. Case  $a, b, c, d = 1, 1, 1, 0$

We first consider scenario 1. For Figure 4, top-left, we fix  $E_2 = E_3 = 1$  and plot the expected utilities of players A and B at the  $\epsilon$ -approximate Nash equilibrium as functions of  $E_1$ , where  $\epsilon = 10^{-3}$  is used. For Figure 4-top-middle, right, we use the same configuration and plot a solution for the probabilities of choosing different resources as a function of  $E_1$  at

the  $\epsilon$ -approximate Nash equilibrium for players A and B, respectively. For scenarios 2 and 3, Figure 4, middle, bottom, have similar descriptions to scenario 1.

We consider the same three scenarios for the simulations on maximizing the worst-case expected utility. In each scenario, for the top figure, we fix  $E_2 = E_3 = 1$  and plot the maximum expected worst-case utility of player A as a function of  $E_1$ . For the bottom figure, we use the same configuration and plot a solution for the probabilities of choosing different resources for player A as a function of  $E_1$ . Notice that the solutions may not be unique, as discussed in Section V-A. Additionally, for Figure 5-top-middle, right, we also indicate the maximum possible error of the solution calculated using the error bound derived in Theorem 4. For scenarios 2 and 3, we have obtained the solutions by averaging over  $10^2$  independent simulations. Further, we have used  $T = 10^5$ ,  $\alpha = 4 \times 10^4$ , and  $V = 2 \times 10^2$ .

Notice that it is difficult to compare the worst-case strategy and the  $\epsilon$ -approximate Nash equilibrium strategy in general since the first can be computed without any cooperation between the players, whereas computing the second requires cooperation among players. Further, as described in Section V-A, the worst-case strategy can be arbitrarily worse than the Nash equilibrium strategy. Nevertheless, comparing Figures 4-left and 5-top, it can be seen that the worst-case strategy and the strategy at  $\epsilon$ -approximate Nash equilibrium yield comparable expected utilities for player A when  $E_1 \geq 2$ . For instance, in scenario 1, for  $E_1 \geq 2$ , the approximate Nash equilibrium strategy coincides with the worst-case strategy of choosing resource 1 with probability 1. However, it should be noted that our algorithm for finding the  $\epsilon$ -approximate Nash equilibrium does not necessarily converge to a socially optimal solution. For instance, in scenario 1, when  $E_1 = 2$ , player A chooses resource 1 with probability 1 and player B chooses resource 2 with probability 1 gives a higher utility for player A without changing the utility of player B.

In Figure 5, it is interesting to notice the variation in choice probabilities of different resources with  $E_1$ . Notice that in scenario 1, the choice probability of resource 1 is non-decreasing for  $E_1 \in [0.1, 0.8]$ , non-increasing for  $E_1 \in [0.8, 1.9]$ , and non-decreasing for  $E_1 \geq 1.9$ . Similar behavior can also be observed for scenario 3. This is surprising since intuition suggests that the probability of choosing a resource should increase with the increasing mean of the reward random variable. However, notice that in scenarios 1 and 3, player B does not observe the reward realization of resource 1. This might force player A, playing for the worst case, to believe that player B increases the probability of choosing resource 1 with increasing  $E_1$ , as a result of which



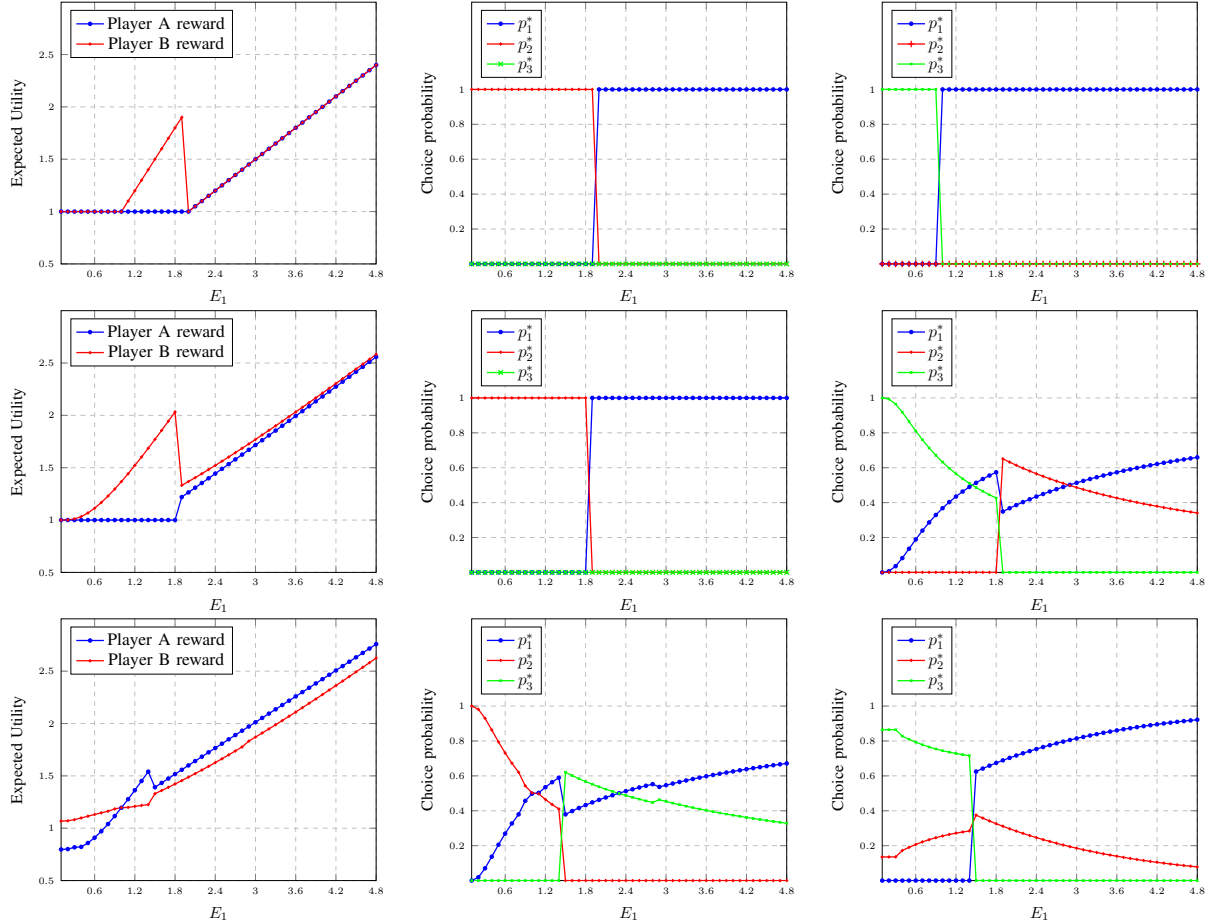


Fig. 4. **Top:** Case  $a = 0, b = 0, c = 3, d = 0$ . **Middle:** Case  $a = 0, b = 1, c = 2, d = 0$ . **Bottom:** Case  $a = b = c = 1, d = 0$ . **Left:** The expected utility of the players at the  $\epsilon$ -approximate Nash equilibrium vs.  $E_1$ . **Middle:** One possible solution for the probabilities of choosing different resources at the  $\epsilon$ -approximate Nash equilibrium for player A vs.  $E_1$ . **Right:** One possible solution for the probabilities of choosing different resources at the  $\epsilon$ -approximate Nash equilibrium for player B vs.  $E_1$ .

player A chooses resource 1 with a lower probability. Notice that the probability of choosing resource 1 in scenario 3 does not grow as fast as the other two. This is because player A observes  $W_1$  and hence can refrain from choosing it when  $W_1$  takes low values.

## VII. CONCLUSIONS

We have implemented the iterative best response algorithm to find the  $\epsilon$ -approximate Nash equilibrium of a two-player stochastic resource-sharing game with asymmetric information. To handle situations where the players do not trust each other and place no assumptions on the incentives of the opponent, we solved the problem of maximizing the worst-case expected

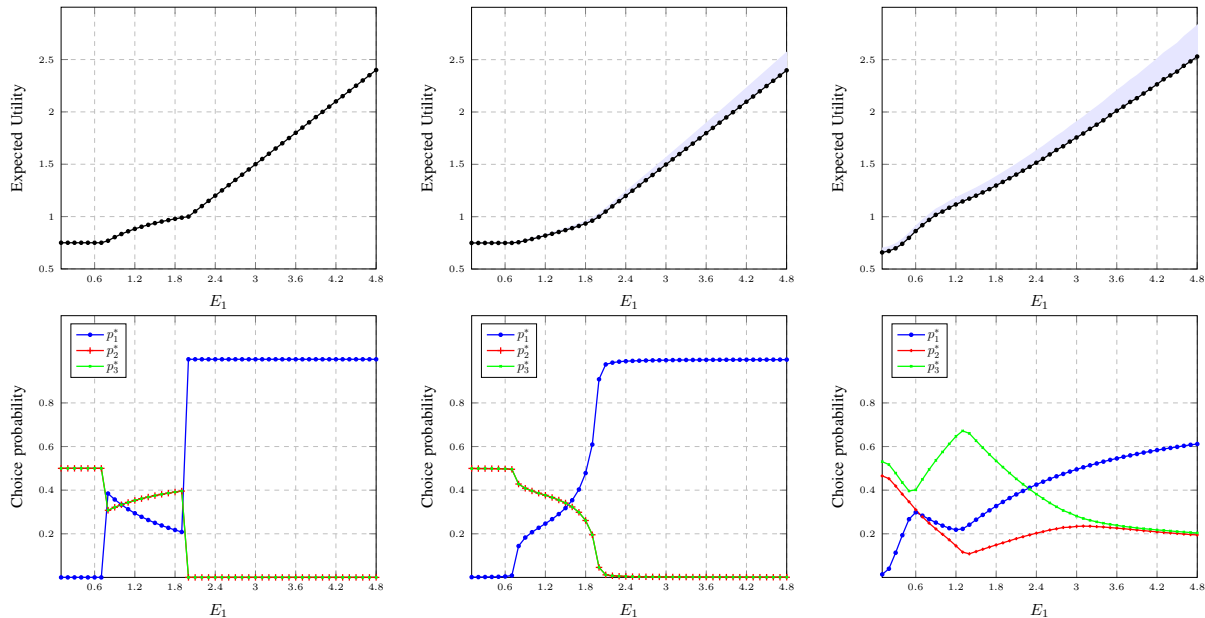


Fig. 5. **Left:** Case  $a = 0, b = 0, c = 3, d = 0$ . **Middle:** Case  $a = 0, b = 1, c = 2, d = 0$ . **Right:** Case  $a = b = c = 1, d = 0$ . **Top:** The maximum expected worst-case utility of player A and the error margin (shaded in blue) vs.  $E_1$ . **Bottom:** One possible solution for the probabilities of choosing different resources for player A vs.  $E_1$ .

utility of the first player using a novel algorithm that combines drift-plus penalty theory and online optimization techniques. An explicit solution can be constructed when both players do not observe the realizations of any of the reward random variables. This special case leads to counter-intuitive insights.

In our approach, we have assumed that the reward random variables of different resources are independent. It should be noted that this assumption can be relaxed without affecting the analysis for the special case when both players do not observe the realizations of any of the reward random variables. An interesting question would be what happens in the general case when the reward random variables are not independent. While it is still possible to implement our algorithm in this setting, it is not guaranteed that the algorithm will converge to the optimal solution. Hence, finding an algorithm for this case that exploits the correlations between the reward random variables could be potential future work.

Several other extensions can be considered as well. One would be considering a scenario with multiple players. The general multiplayer case yields a complex information structure since the set of resources has to be split into  $2^m$  subsets, where  $m$  is the number of players. Additionally,

the idea of conditioning on the common information is difficult to be adapted for this case. Nevertheless, various simplified schemes could be considered. One example would be a case with no common information. In this case, the set of resources is split into  $m + 1$  disjoint subsets where the  $i$ -th ( $1 \leq i \leq m$ ) subset is the subset of resources of which the  $i$ -th player observes the rewards, and the  $m + 1$ -th subset is the subset of resources of which the rewards are observed by none of the players. Another interesting scenario is when no player observes any of the reward realizations. In both these cases, the expected utility can be calculated following a similar procedure to the two-player case, but finding the worst-case expected utility is difficult. Hence, we believe both cases could be potential future work. Another extension would be extending the algorithm to be implemented with a repeated game structure and in an online scenario.

## APPENDIX A

### PROOF OF THEOREM 2

Notice that the term  $\mathbb{E}\{\max\{\Omega_j x_j; 1 \leq j \leq n\} | \mathbf{Z}\}$  of  $f$  is convex since the max function is convex and expectation preserves convexity. Hence  $f$  is concave.

For the 2, and 3, we use the two inequalities,

$$f(\mathbf{x}) - f(\mathbf{y}) \geq \sum_{j \in \mathcal{A}} (x_j - y_j) + \sum_{j \in \mathcal{A}^c} E_j (x_j - y_j) - \frac{1}{2} \mathbb{E}\{\max\{\Omega_j (x_j - y_j); j \in [1 : n]\} | \mathbf{Z}\}, \quad (70)$$

and

$$f(\mathbf{x}) - f(\mathbf{y}) \leq \sum_{j \in \mathcal{A}} (x_j - y_j) + \sum_{j \in \mathcal{A}^c} E_j (x_j - y_j) + \frac{1}{2} \mathbb{E}\{\max\{\Omega_j (y_j - x_j); j \in [1 : n]\} | \mathbf{Z}\}, \quad (71)$$

both of which follow from the fact that for real numbers  $\gamma_1, \gamma_2, \gamma_3, \gamma_4$ ,  $\max\{\gamma_1 + \gamma_2, \gamma_3 + \gamma_4\} \leq \max\{\gamma_1, \gamma_3\} + \max\{\gamma_2, \gamma_4\}$ .

For 2, we consider  $\mathbf{x} \geq \mathbf{y}$  where the inequality is entry-wise. Notice that,

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{y}) &\geq \sum_{j \in \mathcal{A}} (x_j - y_j) + \sum_{j \in \mathcal{A}^c} E_j (x_j - y_j) - \frac{1}{2} \mathbb{E}\left\{ \sum_{j=1}^n \Omega_j (x_j - y_j) \middle| \mathbf{Z} \right\} \\ &= \sum_{j \in \mathcal{A}} \frac{1}{2} (x_j - y_j) + \sum_{j \in \mathcal{A}^c} \frac{E_j}{2} (x_j - y_j), \end{aligned}$$

where the inequality follows from (70) and the fact that for  $\gamma_1, \gamma_2 \geq 0$ ,  $\max\{\gamma_1, \gamma_2\} \leq \gamma_1 + \gamma_2$ .

For 3, note that

$$f(\mathbf{x}) - f(\mathbf{y}) \leq_{(a)} \sum_{j \in \mathcal{A}} |x_j - y_j| + \sum_{j \in \mathcal{A}^c} E_j |x_j - y_j| + \frac{1}{2} \mathbb{E}\left\{ \sum_{j=1}^n \Omega_j |x_j - y_j| \middle| \mathbf{Z} \right\}$$

$$= \frac{3}{2} \sum_{j \in \mathcal{A}} |x_j - y_j| + \frac{3}{2} \sum_{j \in \mathcal{A}^c} E_j |x_j - y_j|, \quad (72)$$

where (a) follows from (71) and the fact that for  $\gamma_1, \gamma_2 \geq 0$ ,  $\max\{\gamma_1, \gamma_2\} \leq \gamma_1 + \gamma_2$ .

## APPENDIX B

### PROOF OF THEOREM 3

Define the vector-valued function  $P^A : \mathbb{R}^{a+d} \mapsto \mathbb{R}^n$  where,

$$P_k^A(\mathbf{x}, \mathbf{z}) = \begin{cases} \Pr\{\alpha^A = k | \mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z}\} & \text{if } \Pr\{\mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z}\} > 0, \\ 0 & \text{otherwise,} \end{cases} \quad (73)$$

for  $k \in [1 : n]$ .

Notice that the function  $P^A$  completely determines the probabilistic decision of player A. We call  $P^A$  the strategy function of player A.

Take any  $(\boldsymbol{\alpha}, \boldsymbol{\psi}), (\boldsymbol{\eta}, \boldsymbol{\theta}) \in \mathcal{G}^A$ . Then there exists strategies  $g^{A,1}$  and  $g^{A,2}$  with strategy functions  $P^{A,1}$ , and  $P^{A,2}$  such that,

$$\begin{aligned} \alpha_k &= \mathbb{E}\{P_k^{A,1}(\mathbf{X}, \mathbf{Z}) | \mathbf{Z}\}, \quad \eta_k = \mathbb{E}\{P_k^{A,2}(\mathbf{X}, \mathbf{Z}) | \mathbf{Z}\} \text{ for } 1 \leq k \leq n, \\ \psi_l &= \mathbb{E}\{W_l P_l^{A,1}(\mathbf{X}, \mathbf{Z}) | \mathbf{Z}\}, \quad \theta_l = \mathbb{E}\{W_l P_l^{A,2}(\mathbf{X}, \mathbf{Z}) | \mathbf{Z}\} \text{ for } 1 \leq l \leq a. \end{aligned} \quad (74)$$

Consider  $\lambda \in [0, 1]$ , and define  $P^{A,*} = \lambda P^{A,1} + (1 - \lambda) P^{A,2}$ . Hence,

$$\lambda(\boldsymbol{\alpha}, \boldsymbol{\psi}) + (1 - \lambda)(\boldsymbol{\eta}, \boldsymbol{\theta}) = (\lambda\boldsymbol{\alpha} + (1 - \lambda)\boldsymbol{\eta}, \lambda\boldsymbol{\psi} + (1 - \lambda)\boldsymbol{\theta}) = (\boldsymbol{\beta}, \boldsymbol{\gamma}), \quad (75)$$

where,

$$\begin{aligned} \beta_k &= \mathbb{E}\{P_k^{A,*}(\mathbf{X}, \mathbf{Z}) | \mathbf{Z}\} \text{ for } 1 \leq k \leq n, \\ \gamma_l &= \mathbb{E}\{W_l P_l^{A,*}(\mathbf{X}, \mathbf{Z}) | \mathbf{Z}\} \text{ for } 1 \leq l \leq a. \end{aligned} \quad (76)$$

Notice that  $P^{A,*}$  is the strategy function of the strategy  $g^{A,*}$  of using  $g^{A,1}$  and  $g^{A,2}$  in a mixture with probabilities  $\lambda$  and  $1 - \lambda$ . Hence  $(\boldsymbol{\beta}, \boldsymbol{\gamma}) \in \mathcal{G}^A$ , which completes the proof.

## APPENDIX C

## PROOF OF LEMMA 3

We begin with several results which are used in the proof.

*Lemma 10:* If  $(\mathbf{p}^*, \gamma^*)$  solves the problem,

$$\begin{aligned} \text{(P2-1): maximize} \quad & \sum_{k=1}^n p_k E_k - \frac{1}{2}\gamma \\ \text{subject to} \quad & \mathbf{p} \in \mathcal{I}, \\ & \gamma \geq p_k E_k \quad \forall 1 \leq k \leq n, \end{aligned} \tag{77}$$

where  $\mathcal{I}$  is the  $n$ -dimensional probability simplex defined in (32), then  $\mathbf{p}^*$  solves (P2).

*Proof of Lemma 10:* Define,

$$f_1(\mathbf{p}, \gamma) = \sum_{k=1}^n p_k E_k - \frac{1}{2}\gamma. \tag{78}$$

Notice that  $f(\mathbf{p}) = f_1(\mathbf{p}, \max\{p_k E_k; 1 \leq k \leq n\})$ . Let  $(\mathbf{p}^*, \gamma^*)$  be a solution for (P2-1). Notice that for  $(\mathbf{p}^*, \gamma^*)$  to be feasible for (P2-1), we should have  $\gamma^* \geq \max\{p_k^* E_k; 1 \leq k \leq n\}$ . However, if  $\gamma^* > \max\{p_k^* E_k; 1 \leq k \leq n\}$ , we have that  $f_1(\mathbf{p}^*, \max\{p_k^* E_k; 1 \leq k \leq n\}) > f_1(\mathbf{p}^*, \gamma^*)$ , which contradicts the optimality of  $(\mathbf{p}^*, \gamma^*)$  for (P2-1). Hence,  $\gamma^* = \max\{p_k^* E_k; 1 \leq k \leq n\}$ . Hence, we have  $f(\mathbf{p}^*) = f_1(\mathbf{p}^*, \gamma^*)$ .

Now, consider  $\tilde{\mathbf{p}} \in \mathcal{I}$ . Define  $\tilde{\gamma} = \max\{\tilde{p}_k E_k; 1 \leq k \leq n\}$ . Since  $(\tilde{\mathbf{p}}, \tilde{\gamma})$  is also feasible for (P2-1), we should have  $f_1(\tilde{\mathbf{p}}, \tilde{\gamma}) \leq f_1(\mathbf{p}^*, \gamma^*)$ . This implies  $f(\mathbf{p}^*) \geq f(\tilde{\mathbf{p}})$ . Hence,  $\mathbf{p}^*$  is an optimal solution of (P2). ■

*Lemma 11:* Consider fixed  $\boldsymbol{\mu} \in \mathbb{R}^n$  such that  $\mu_k \geq 0$  for all  $1 \leq k \leq n$ . Now, consider the unconstrained problem,

$$\begin{aligned} \text{(P2-2): maximize} \quad & f_2(\mathbf{p}, \gamma) = \sum_{k=1}^n p_k E_k - \frac{1}{2}\gamma + \sum_{k=1}^n \mu_k (\gamma - p_k E_k) \\ \text{subject to} \quad & \mathbf{p} \in \mathcal{I}, \gamma \in \mathbb{R}. \end{aligned} \tag{79}$$

Assume  $(\mathbf{p}^*, \gamma^*)$  is a solution (P2-2). Additionally, assume that

$$\begin{aligned} E_k p_k^* &\leq \gamma^* \quad \text{for all } 1 \leq k \leq n, \\ E_k p_k^* &= \gamma^* \quad \text{whenever } \mu_k > 0. \end{aligned} \tag{80}$$

Then  $(\mathbf{p}^*, \gamma^*)$  is a solution for (P2-1).

*Proof of Lemma 11:* First, notice that  $(\mathbf{p}^*, \gamma^*)$  satisfies the constraints of (P2-1). To show that it maximizes the objective in (P2-1), consider any  $(\mathbf{p}, \gamma)$  that is feasible for (P2-1). Notice that

$$\begin{aligned}
f_1(\mathbf{p}, \gamma) &= f_2(\mathbf{p}, \gamma) - \sum_{\mu_k > 0} \mu_k (\gamma - p_k E_k) \\
&\stackrel{(a)}{\leq} f_2(\mathbf{p}^*, \gamma^*) - \sum_{\mu_k > 0} \mu_k (\gamma - p_k E_k) \\
&= f_1(\mathbf{p}^*, \gamma^*) + \sum_{\mu_k > 0} \mu_k (\gamma^* - p_k^* E_k - \gamma + p_k E_k) \\
&\stackrel{(b)}{=} f_1(\mathbf{p}^*, \gamma^*) - \sum_{\mu_k > 0} \mu_k (\gamma - p_k E_k) \stackrel{(c)}{\leq} f_1(\mathbf{p}^*, \gamma^*), \tag{81}
\end{aligned}$$

where  $f_1$  is the objective of (P2-1) defined in (78),  $f_2$  is the objective of (P2-2), (a) follows from the optimality of  $(\mathbf{p}^*, \gamma^*)$  for (P2-2), (b) follows due to (80), and (c) follows since  $\mu_k \geq 0$  and  $(\mathbf{p}, \gamma)$  is feasible for (P2-1). Hence, we have the result.  $\blacksquare$

Define,

$$S_k = \sum_{j=1}^k \frac{1}{E_j}, \tag{82}$$

for  $1 \leq k \leq n$ . We also establish the following lemma, which is useful in our solution.

*Lemma 12:* Let

$$r = \arg \max_{1 \leq k \leq n} \frac{k - \frac{1}{2}}{S_k}, \tag{83}$$

where  $\arg \max$  returns the lowest index in the case of ties. Let us also define  $\boldsymbol{\mu} \in \mathbb{R}^n$  as

$$\mu_k = \begin{cases} 1 - \frac{1}{E_k} \frac{r - \frac{1}{2}}{S_r} & \text{if } 1 \leq k \leq r, \\ 0 & \text{otherwise.} \end{cases} \tag{84}$$

Then we have

- 1)  $\mu_k \geq 0$  for all  $k$  such that  $1 \leq k \leq n$ .
- 2)  $\sum_{k=1}^n \mu_k = \frac{1}{2}$ .
- 3)  $E_k(1 - \mu_k) = \frac{r - \frac{1}{2}}{S_r}$  for  $1 \leq k \leq r$ .
- 4)  $E_k(1 - \mu_k) \leq \frac{r - \frac{1}{2}}{S_r}$  for  $r + 1 \leq k \leq n$ .

*Proof of Lemma 12:*

- 1) Notice that by the definition of  $\mu_k$ , it is enough to prove the result for  $1 \leq k \leq r$ . Notice that we are required to prove that

$$\frac{1}{E_k} \frac{r - \frac{1}{2}}{S_r} \leq 1, \quad (85)$$

for all  $1 \leq k \leq r$ . Since  $E_k \geq E_{k+1}$  for  $1 \leq k \leq n - 1$ , it suffices to prove that

$$\frac{1}{E_r} \frac{r - \frac{1}{2}}{S_r} \leq 1. \quad (86)$$

We consider two cases.

**Case 1:**  $r = 1$ . This case reduces to,

$$\frac{1}{2E_1} \leq \frac{1}{E_1}, \quad (87)$$

which is trivial.

**Case 2:**  $r > 1$ . Note that from the definition of  $r$  in (83), we have

$$\frac{r - \frac{1}{2}}{S_r} \geq \frac{r - \frac{3}{2}}{S_{r-1}}. \quad (88)$$

After substituting  $S_{r-1} = S_r - \frac{1}{E_r}$  and rearranging, we have the desired result.

- 2) Notice that

$$\sum_{k=1}^n \mu_k = \sum_{k=1}^r \mu_k = \sum_{k=1}^r \left( 1 - \frac{1}{E_k} \frac{r - \frac{1}{2}}{S_r} \right) = r - \frac{r - \frac{1}{2}}{S_r} \sum_{k=1}^r \frac{1}{E_k} = r - \frac{r - \frac{1}{2}}{S_r} S_r = \frac{1}{2}. \quad (89)$$

- 3) This follows from the definition of  $\mu_k$  for  $1 \leq k \leq r$ .
- 4) There is nothing to prove if  $r = n$ . Hence, we can assume  $r < n$ . Since  $\mu_k = 0$  for  $k \geq r + 1$ , it suffices to prove that  $E_k \leq \frac{r - (1/2)}{S_r}$ . Notice that if we can prove the result for  $k = r + 1$ , we are finished since  $E_k \geq E_{k+1}$  for  $1 \leq k \leq n$ . Note that from the definition of  $r$  in (83), we have

$$\frac{r - \frac{1}{2}}{S_r} \geq \frac{r + \frac{1}{2}}{S_{r+1}}. \quad (90)$$

After substituting  $S_{r+1} = S_r + \frac{1}{E_{r+1}}$  and rearranging, we have the desired result. ■

Now, we solve the problem using the above lemmas. Consider the problem defined in Lemma 11 with  $\mu$  defined in Lemma 12. Specifically, consider the problem,

$$\begin{aligned} \text{(P2-3): maximize } \quad & f_2(\mathbf{p}, \gamma) = \sum_{k=1}^n p_k E_k - \frac{1}{2} \gamma + \sum_{k=1}^n \mu_k (\gamma - p_k E_k) \\ \text{subject to } \quad & \mathbf{p} \in \mathcal{I}, \gamma \in \mathbb{R}. \end{aligned} \quad (91)$$

where  $\boldsymbol{\mu}$  and  $r$  are defined in (83)-(84). For this choice of  $\mu_k$  we have

$$f_2(\boldsymbol{p}, \gamma) = \sum_{k=1}^n p_k E_k (1 - \mu_k) + \gamma \left( \sum_{k=1}^n \mu_k - \frac{1}{2} \right) = \sum_{k=1}^n p_k E_k (1 - \mu_k), \quad (92)$$

where the last equality follows from Lemma 12-2. Now, due to Lemma 12-3 and Lemma 12-4, the optimal solution for (P2-3) is any  $(\boldsymbol{p}, \gamma)$  such that  $\gamma \in \mathbb{R}$ , and  $\boldsymbol{p} \in \mathcal{I}$  such that  $p_k = 0$  for  $k > r$ . In particular, consider the solution  $(\boldsymbol{p}^*, \gamma^*)$  given by,

$$p_k^* = \begin{cases} \frac{1}{E_k S_r} & \text{if } k \leq r, \\ 0 & \text{otherwise,} \end{cases} \quad (93)$$

and  $\gamma^* = \frac{1}{S_r}$ . Notice that for  $1 \leq k \leq r$ , we have that  $p_k^* E_k = \gamma^*$ , and  $p_k^* E_k = 0 \leq \gamma^*$  for  $r+1 \leq k \leq n$ . Hence, from Lemma 11,  $(\boldsymbol{p}^*, \gamma^*)$  is a solution for (P2-1). Hence, from Lemma 10,  $\boldsymbol{p}^*$  is a solution for (P2) as desired.

#### APPENDIX D

##### PROOF OF LEMMA 4

Notice that,

$$\begin{aligned} \Delta(t) &= \mathbb{E}\{L(t+1) - L(t) | \mathcal{H}(t), \mathbf{Z}\} = \frac{1}{2} \mathbb{E} \left\{ \sum_{j=1}^n Q_j(t+1)^2 - Q_j(t)^2 \middle| \mathcal{H}(t), \mathbf{Z} \right\} \\ &= \frac{1}{2} \sum_{j=1}^n \mathbb{E}\{Q_j(t+1)^2 | \mathcal{H}(t), \mathbf{Z}\} - \frac{1}{2} \sum_{j=1}^n Q_j(t)^2 \\ &= \frac{1}{2} \sum_{j=1}^n \mathbb{E}\{\max\{Q_j(t) + \gamma_j(t) - x_j(t), 0\}^2 | \mathcal{H}(t), \mathbf{Z}\} - \frac{1}{2} \sum_{j=1}^n Q_j(t)^2 \\ &\leq \frac{1}{2} \sum_{j=1}^n \mathbb{E}\{(Q_j(t) + \gamma_j(t) - x_j(t))^2 | \mathcal{H}(t), \mathbf{Z}\} - \frac{1}{2} \sum_{j=1}^n Q_j(t)^2 \\ &\leq \sum_{j=1}^n Q_j(t) (\gamma_j(t) - \mathbb{E}\{x_j(t) | \mathcal{H}(t), \mathbf{Z}\}) + \frac{1}{2} \sum_{j=1}^n \gamma_j(t)^2 + \frac{1}{2} \sum_{j=1}^n \mathbb{E}\{x_j(t)^2 | \mathcal{H}(t), \mathbf{Z}\} \\ &\stackrel{(a)}{\leq} \sum_{j=1}^n Q_j(t) (\gamma_j(t) - \mathbb{E}\{x_j(t) | \mathcal{H}(t), \mathbf{Z}\}) + \frac{1}{2} \sum_{j \in \mathcal{A}} E_j^2 + \frac{n-a}{2} \\ &\quad + \frac{1}{2} \sum_{j \in \mathcal{A}} \mathbb{E}\{(X_j(t) \mathbb{1}_{\{\alpha(t)=k\}})^2 | \mathcal{H}(t), \mathbf{Z}\} + \frac{1}{2} \sum_{j \in \mathcal{A}^c} \mathbb{E}\{(\mathbb{1}_{\{\alpha(t)=k\}})^2 | \mathcal{H}(t), \mathbf{Z}\} \\ &\leq \sum_{j=1}^n Q_j(t) (\gamma_j(t) - \mathbb{E}\{x_j(t) | \mathcal{H}(t), \mathbf{Z}\}) + \frac{1}{2} \sum_{j \in \mathcal{A}} E_j^2 + \frac{n-a}{2} + \frac{1}{2} \sum_{j \in \mathcal{A}} \mathbb{E}\{X_j(t)^2 | \mathcal{H}(t), \mathbf{Z}\} \end{aligned}$$



$$\begin{aligned}
& + \frac{1}{2} \sum_{j \in \mathcal{A}^c} \mathbb{E}\{1 | \mathcal{H}(t), \mathbf{Z}\} \\
& \stackrel{(b)}{=} \sum_{j=1}^n Q_j(t) (\gamma_j(t) - \mathbb{E}\{x_j(t) | \mathcal{H}(t), \mathbf{Z}\}) + \frac{1}{2} \sum_{j \in \mathcal{A}} E_j^2 + \frac{n-a}{2} + \frac{1}{2} \sum_{j \in \mathcal{A}} \mathbb{E}\{W_j^2\} + \frac{n-a}{2} \\
& = \sum_{j=1}^n Q_j(t) (\gamma_j(t) - \mathbb{E}\{x_j(t) | \mathcal{H}(t), \mathbf{Z}\}) + D_1, \tag{94}
\end{aligned}$$

where the inequality (a) follows since  $\gamma(t) \in \mathcal{K}$  and equality (b) follows from the fact that  $\mathbf{X}(t)$  is independent of  $\mathcal{H}(t)$  and  $\mathbf{Z}$ .

## APPENDIX E

### PROOF OF LEMMA 6

Notice that,

$$f'_t(\gamma(t-1)) = \mathbf{v} - \frac{1}{2} \tilde{\Omega}(t), \tag{95}$$

where  $\mathbf{v}$  is defined by,

$$v_j = \begin{cases} 1 & \text{if } j \in \mathcal{A}, \\ E_j & \text{if } j \in \mathcal{A}^c. \end{cases} \tag{96}$$

$\tilde{\Omega}(t)$  is given by,  $\tilde{\Omega}_k(t) = \Omega_k(t) \mathbb{1}_{\{\arg \max_{1 \leq j \leq n} \{\gamma_j(t-1)\Omega_j(t)\} = k\}}$ , and  $\arg \max$  return the least index in case of ties. Notice that

$$\begin{aligned}
& -V f'_t(\gamma(t-1))^\top \{\gamma(t) - \gamma(t-1)\} + \alpha \|\gamma(t) - \gamma(t-1)\|_2^2 \\
& \geq_{(a)} -V \|f'_t(\gamma(t-1))\|_2 \|\gamma(t) - \gamma(t-1)\|_2 + \alpha \|\gamma(t) - \gamma(t-1)\|_2^2 \\
& = \alpha \left( \|\gamma(t) - \gamma(t-1)\|_2 - \frac{V}{2\alpha} \|f'_t(\gamma(t-1))\|_2 \right)^2 - \frac{V^2}{4\alpha} \|f'_t(\gamma(t-1))\|_2^2 \\
& \geq -\frac{V^2}{4\alpha} \|f'_t(\gamma(t-1))\|_2^2 = -\frac{V^2}{4\alpha} \left\| \mathbf{v} - \frac{1}{2} \tilde{\Omega}(t) \right\|_2^2 \geq_{(b)} -\frac{V^2}{4\alpha} \|\mathbf{v}\|_2^2 - \frac{V^2}{16\alpha} \|\tilde{\Omega}(t)\|_2^2 \\
& \geq -\frac{V^2}{4\alpha} \left( a + \sum_{j \in \mathcal{A}^c} E_j^2 \right) - \frac{V^2}{16\alpha} \|\Omega(t)\|_2^2, \tag{97}
\end{aligned}$$

where (a) follows from the Cauchy-Schwarz inequality, (b) follows since  $v_k \geq 0$ , and  $\tilde{\Omega}_k(t) \geq 0$  for all  $1 \leq k \leq n$ .

APPENDIX F  
PROOF OF LEMMA 8

Notice that from the definition of  $Q_j(t+1)$  in (39) and the definition of  $x_j(t)$  in (46) we have that,

$$Q_j(t+1) - Q_j(t) \geq \gamma_j(t) - x_j(t), \quad (98)$$

for all  $1 \leq j \leq n$  and  $1 \leq t \leq T-1$ . Summing the above from 1 to  $T-1$  we have that,

$$Q_j(T) - Q_j(1) \geq \sum_{t=1}^{T-1} \{\gamma_j(t) - x_j(t)\}. \quad (99)$$

After using  $Q_j(1) = 0$ , taking expectations conditioned on  $\mathbf{Z}$ , and some algebraic manipulations we have,

$$\frac{\mathbb{E}\{Q_j(T)|\mathbf{Z}\}}{T} \geq \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\{\gamma_j(t) - x_j(t)|\mathbf{Z}\}. \quad (100)$$

We have the desired inequality from the above since  $Q_j(T)$  is non-negative.

APPENDIX G  
PROOF OF LEMMA 9

Define the  $\mathbf{v}, \mathbf{u}$  as follows.

$$v_k = \begin{cases} 1 & \text{if } k \in \mathcal{A}, \\ E_k & \text{if } k \in \mathcal{A}^c, \end{cases} \quad \text{and } u_k = \begin{cases} E_k & \text{if } k \in \mathcal{A}, \\ 1 & \text{if } k \in \mathcal{A}^c. \end{cases} \quad (101)$$

Hence we are required to prove that  $Q_j(t) \leq (v_j + 2\sqrt{2}u_j)\sqrt{\alpha} + u_j$  for all  $t \in [1 : T]$ .

We begin with several important results.

*Lemma 13:* We have the following results regarding  $Q_j(t)$ .

- 1)  $Q_j(t+1) \leq Q_j(t) + u_j$  for all  $t \geq 1$ .
- 2) Assume  $Q_j(t) \geq (v_j + \sqrt{2}u_j)\sqrt{\alpha}$  for some  $t \geq 1$ . Then we have, either  $\gamma_j(t) = 0$  or

$$\gamma_j(t) \leq \gamma_j(t-1) - \frac{u_j}{\sqrt{2\alpha}}. \quad (102)$$

- 3) Assume  $Q_j(\tau) \geq (v_j + \sqrt{2}u_j)\sqrt{\alpha}$  for all  $\tau \in [t : t + t_0]$ , where  $t \geq 1$  and  $t_0 \geq 0$ . Additionally assume  $\gamma_j(t-1) = 0$ . Then  $\gamma_j(\tau) = 0$  for all  $\tau \in [t-1 : t + t_0]$ .

*Proof:*

1) Notice that from the definition of  $Q_j(t+1)$  in (39), for  $j \in \mathcal{A}$  we have

$$\begin{aligned} Q_j(t+1) &= \max \{Q_j(t) + \gamma_j(t) - X_j(t) \mathbb{1}_{\{\alpha^A(t)=j\}}, 0\} \leq \max \{Q_j(t) + u_j, 0\} \\ &= Q_j(t) + u_j, \end{aligned} \quad (103)$$

where the inequality follows from the definition of  $u_j$  in (101). The same argument can be repeated for  $j \in \mathcal{A}^c$ .

2) Notice that if  $\gamma_j(t) \neq 0$  then we have

$$\gamma_j(t) \leq \gamma_j(t-1) - \frac{-V f'_{t,j}(\gamma(t-1)) + Q_j(t)}{2\alpha}, \quad (104)$$

which follows since  $\gamma_j(t)$  is the projection of  $\gamma_j(t-1) - \frac{-V f'_{t,j}(\gamma(t-1)) + Q_j(t)}{2\alpha}$  onto  $[0, u_j]$  (See (42)). Hence we have that,

$$\begin{aligned} \gamma_j(t) &\leq \gamma_j(t-1) - \frac{-V f'_{t,j}(\gamma(t-1)) + Q_j(t)}{2\alpha} \\ &\leq_{(a)} \gamma_j(t-1) - \frac{-V v_j + (v_j + \sqrt{2}u_j)\sqrt{\alpha}}{2\alpha} \leq_{(b)} \gamma_j(t-1) - \frac{u_j}{\sqrt{2\alpha}}, \end{aligned} \quad (105)$$

where (a) follows from the subgradients of  $f_t$  found in (38) and (b) follows from  $\alpha \geq V^2$ .

3) Notice if we prove  $\gamma_j(t) = 0$ , we can use the same argument inductively to establish the result. Assume the contrary that  $\gamma_j(t) \neq 0$ . Then from part-2, we should have,

$$\gamma_j(t) \leq \gamma_j(t-1) - \frac{u_j}{\sqrt{2\alpha}} = -\frac{u_j}{\sqrt{2\alpha}}, \quad (106)$$

which is a contradiction since  $\gamma_j(t) \geq 0$ . Hence we have the result. ■

Now we use an inductive argument to prove the main result. Notice that the result is true for  $t = 1$ , since  $Q_j(1) = 0 \leq (v_j + 2\sqrt{2}u_j)\sqrt{\alpha} + u_j$ . Now we prove that  $Q_j(t+1) \leq (v_j + 2\sqrt{2}u_j)\sqrt{\alpha} + u_j$  for  $t \geq 1$ , with the assumption that,  $Q_j(t) \leq (v_j + 2\sqrt{2}u_j)\sqrt{\alpha} + u_j$ .

We consider three cases.

**Case 1:**  $Q_j(t) \leq (v_j + 2\sqrt{2}u_j)\sqrt{\alpha}$ . This case follows from Lemma 13-1.

**Case 2:**  $t \leq \sqrt{2\alpha} + 1$ . Notice that,

$$Q_j(t+1) \leq Q_j(1) + u_j t \leq (\sqrt{2\alpha} + 1)u_j \leq (v_j + 2\sqrt{2}u_j)\sqrt{\alpha} + u_j, \quad (107)$$

where first inequality follows from Lemma 13-1.

**Case 3:**  $t > \sqrt{2\alpha} + 1$  and  $Q_j(t) > (v_j + 2\sqrt{2}u_j)\sqrt{\alpha}$ . For this, we prove that  $\gamma_j(t) = 0$ , which establishes the claim from the definition of  $Q_j(t+1)$  in (39) and the induction hypothesis.

Notice that for all  $u \in [1 : t]$  we have,

$$\begin{aligned} Q_j(u) &\geq_{(a)} Q_j(t) - (t-u)u_j \geq (v_j + 2\sqrt{2}u_j)\sqrt{\alpha} - (t-u)u_j \\ &= (v_j + \sqrt{2}u_j)\sqrt{\alpha} + \sqrt{2\alpha}u_j - (t-u)u_j, \end{aligned} \quad (108)$$

where (a) follows from Lemma 13-1.

Hence, for all  $u \in \mathbb{Z}$  such that  $t - \sqrt{2\alpha} \leq u \leq t$ , we have that,

$$Q_j(u) \geq (v_j + \sqrt{2}u_j)\sqrt{\alpha}. \quad (109)$$

Now we prove that there exists  $u \in \mathbb{Z}$  such that  $t - \sqrt{2\alpha} \leq u \leq t$  and  $\gamma_j(u) = 0$ , which will establish that  $\gamma_j(t) = 0$  from Lemma 13-3. For the proof assume the contrary  $\gamma_j(u) > 0$  for all  $u \in \mathbb{Z}$  such that  $t - \sqrt{2\alpha} \leq u \leq t$  (or  $u \in [t - \lfloor \sqrt{2\alpha} \rfloor : t]$ , where  $\lfloor x \rfloor$  denotes the largest integer smaller than or equal to  $x$ ). From Lemma 13-2 we have that,

$$\gamma_j(t) \leq \gamma_j\left(t - \lfloor \sqrt{2\alpha} \rfloor - 1\right) - \frac{(\lfloor \sqrt{2\alpha} \rfloor + 1)u_j}{\sqrt{2\alpha}} \leq 0, \quad (110)$$

where the last inequality follows since  $\lfloor x \rfloor + 1 \geq x$  and  $\gamma_j(t - \lfloor \sqrt{2\alpha} \rfloor - 1) \leq u_j$  (since  $\gamma_j(\tau) \in [0, u_j]$  for all  $\tau \in [1 : T]$  by the projection definition of  $\gamma_j(\tau)$  in (42), and  $\gamma_j(0) = 0$ ). Hence, we should have that  $\gamma_j(t) = 0$  which contradicts our initial assumption. Hence, we are done.

## APPENDIX H

### SOME SPECIAL CASES IN WHICH SIMPLER ALTERNATIVE SOLUTIONS EXIST

Below, we discuss two special cases in which a simpler solution can be obtained for the worst-case expected utility maximization problem. First, we discuss the case  $a = 0$  (with no restriction on  $b \in \{0, 1, \dots, n\}$ ), after which we discuss the case  $a = 1$ , where  $W_1$  is assumed to have a continuous CDF.

#### A. Mirror descent solution for $a = 0$

Observe that when  $a = 0$  (with no restriction on  $b \in \{0, 1, \dots, n\}$ ), (P1.1) can be written as,

$$\begin{aligned} \text{(P4): maximize } & \sum_{k=1}^n p_k E_k - \frac{1}{2} \mathbb{E}\{\max\{p_k \Omega_k; 1 \leq k \leq n\}\} \\ \text{subject to } & \mathbf{p} \in \mathcal{I}, \end{aligned} \quad (111)$$

where  $\mathcal{I}$  is the simplex set defined in Theorem 3. Define  $\mathcal{I}^0 = \mathcal{I} \cap \mathbb{R}_+^n$ . We use the online mirror descent algorithm to solve this problem. We first formulate the above problem as an online convex optimization problem, where we obtain the solution by iteratively updating  $\mathbf{p}$  for  $T$  iterations. In the  $t$ -th iteration, we sample  $\Omega(t)$  from the distribution of  $\Omega$ , where  $\Omega$  is defined in (21). Let,

$$f_t(\mathbf{p}) = -\sum_{k=1}^n p_k E_k + \frac{1}{2} \max\{p_k \Omega_k(t); 1 \leq k \leq n\}. \quad (112)$$

Let  $\mathbf{p}(t)$  denote the value of  $\mathbf{p}$  after  $t$ -th iteration. we begin with  $\mathbf{p}(0) = [1/n, 1/n, \dots, 1/n]$ , and we obtain  $\mathbf{p}(t)$  ( $1 \leq t \leq T$ ) by solving,

$$\text{(P4-1): minimize}_{\mathbf{p}(t)} \quad \nabla f_t(\mathbf{p}(t-1))^\top (\mathbf{p}(t) - \mathbf{p}(t-1)) + \alpha D(\mathbf{p}(t) \| \mathbf{p}(t-1)) \quad (113)$$

$$\text{subject to } \mathbf{p}(t) \in \mathcal{I},$$

where  $D(\mathbf{x} \| \mathbf{y})$  denotes the Kullback-Leibler divergence between  $\mathbf{x}$  and  $\mathbf{y}$  given by,

$$D(\mathbf{x} \| \mathbf{y}) = \sum_{k=1}^n x_k \ln \left( \frac{x_k}{y_k} \right), \quad (114)$$

for  $\mathbf{x} \in \mathcal{I}$  and  $\mathbf{y} \in \mathcal{I}^0$ ,  $\nabla f_t(\mathbf{x}) \in \mathbb{R}^n$  is the subgradient of  $f_t$  at  $\mathbf{x}$  which is given by,

$$(\nabla f_t(\mathbf{x}))_j = -E_j + \frac{1}{2} \mathbb{1}_{(\arg \max_{1 \leq k \leq n} \{x_k \Omega_k(t)\} = j)} \Omega_j(t), \quad (115)$$

for each  $1 \leq j \leq n$ , and  $\alpha$  is the constant step size. It is assumed that  $\arg \max$  returns the lowest index in the case of ties. It should be noted that  $\mathbf{p}(t)$  can be found explicitly [38]. The solution is given by,

$$p_k(t) = \frac{p_k(t-1) \exp \left( -\frac{(\nabla f_t(\mathbf{p}(t-1)))_k}{\alpha} \right)}{\sum_{j=1}^n p_j(t-1) \exp \left( -\frac{(\nabla f_t(\mathbf{p}(t-1)))_j}{\alpha} \right)}, \quad (116)$$

for each  $k$  such that  $1 \leq k \leq n$ . Note that  $\mathbf{p}(t) \in \mathcal{I}^0$ . This is useful when establishing the performance bound of the algorithm.

We summarize the above algorithm under Algorithm 2 for brevity. We have the following lemma regarding the performance of the above algorithm.

*Lemma 14:* Let  $\mathbf{p} = \frac{1}{T} \sum_{t=1}^T \mathbf{p}(t-1)$ . Then we have that,

$$\mathbb{E}\{f(\mathbf{p}) | \mathbf{Z}\} - f^* \geq -\frac{2\|\mathbf{E}\|_\infty^2 + \frac{1}{2}\mathbb{E}\{\|\Omega\|_\infty^2 | \mathbf{Z}\}}{2\alpha} - \frac{\alpha}{T} \ln(n), \quad (117)$$

---

**Algorithm 2:** Mirror descent algorithm to solve the case  $a = 0$

---

- 1 Initialize  $\mathbf{p}(0) = [1/n, 1/n, \dots, 1/n]$
  - 2 **for** each iteration  $1 \leq t \leq T$  **do**
  - 3     Sample  $\Omega(t)$  from the distribution of  $\Omega$
  - 4     Calculate the subgradient  $\nabla f_t(\mathbf{p}(t-1))$  according to (115)
  - 5     Obtain  $\mathbf{p}(t)$  according to (116)
  - 6 **end**
  - 7 Calculate  $\mathbf{p}^* = \frac{1}{T} \sum_{t=1}^T \mathbf{p}(t-1)$
- 

where  $f^*$  denotes the optimal value of (P4). Hence, given fixed  $\varepsilon > 0$ , for  $\alpha = 1/\varepsilon$ , and  $T \geq 1/\varepsilon^2$ , the error is  $\mathcal{O}(\varepsilon)$ .

*Proof:* We begin with a few lemmas, which are useful in the proof.

*Lemma 15:* We have that,

$$\nabla f_t(\mathbf{p}(t-1))^\top (\mathbf{p}(t) - \mathbf{p}(t-1)) + \alpha D(\mathbf{p}(t) \| \mathbf{p}(t-1)) \geq -\frac{1}{2\alpha} \|\nabla f_t(\mathbf{p}(t-1))\|_\infty^2. \quad (118)$$

*Proof:* Notice that,

$$\begin{aligned} & \nabla f_t(\mathbf{p}(t-1))^\top (\mathbf{p}(t) - \mathbf{p}(t-1)) + \alpha D(\mathbf{p}(t) \| \mathbf{p}(t-1)) \\ & \geq_{(a)} -\|\nabla f_t(\mathbf{p}(t-1))\|_\infty \|\mathbf{p}(t) - \mathbf{p}(t-1)\|_1 + \alpha D(\mathbf{p}(t) \| \mathbf{p}(t-1)) \\ & \geq_{(b)} -\|\nabla f_t(\mathbf{p}(t-1))\|_\infty \|\mathbf{p}(t) - \mathbf{p}(t-1)\|_1 + \frac{1}{2}\alpha \|\mathbf{p}(t) - \mathbf{p}(t-1)\|_1^2 \\ & = \frac{\alpha}{2} \left( \|\mathbf{p}(t) - \mathbf{p}(t-1)\|_1 - \frac{\|\nabla f_t(\mathbf{p}(t-1))\|_\infty}{\alpha} \right)^2 - \frac{\|\nabla f_t(\mathbf{p}(t-1))\|_\infty^2}{2\alpha} \\ & \geq -\frac{1}{2\alpha} \|\nabla f_t(\mathbf{p}(t-1))\|_\infty^2, \end{aligned} \quad (119)$$

where (a) follows from the Cauchy-Schwarz inequality, and (b) follows from Pinsker's inequality. ■

*Lemma 16:* We have,

$$\mathbb{E}\{\|\nabla f_t(\mathbf{p}(t-1))\|_\infty^2 | \mathbf{Z}\} \leq 2\|\mathbf{E}\|_\infty^2 + \frac{1}{2}\mathbb{E}\{\|\Omega\|_\infty^2 | \mathbf{Z}\}. \quad (120)$$

(Notice that the upper bound is finite since  $W_k$  for each  $k \in \mathcal{B}$  has finite variance).

*Proof:* Notice that,

$$\nabla f_t(\mathbf{p}(t-1)) = -\mathbf{E} + \frac{1}{2}\tilde{\Omega}(t), \quad (121)$$

where  $\tilde{\Omega}(t)$  is given by,  $\tilde{\Omega}_k(t) = \Omega_k(t)\mathbb{1}_{(\arg \max_{1 \leq j \leq n} \{x_j \Omega_j(t)\} = k)}$ . Notice that,

$$\begin{aligned} \mathbb{E}\{\|\nabla f_t(\mathbf{p}(t-1))\|_\infty^2 | \mathbf{Z}\} &\leq \mathbb{E}\left\{\left\|-\mathbf{E} + \frac{1}{2}\tilde{\Omega}(t)\right\|_\infty^2 \middle| \mathbf{Z}\right\} \\ &\leq_{(a)} 2\|\mathbf{E}\|_\infty^2 + \frac{1}{2}\mathbb{E}\left\{\left\|\tilde{\Omega}(t)\right\|_\infty^2 \middle| \mathbf{Z}\right\} \leq 2\|\mathbf{E}\|_\infty^2 + \frac{1}{2}\mathbb{E}\left\{\|\Omega(t)\|_\infty^2 \middle| \mathbf{Z}\right\} \\ &= 2\|\mathbf{E}\|_\infty^2 + \frac{1}{2}\mathbb{E}\{\|\Omega\|_\infty^2 | \mathbf{Z}\}, \end{aligned} \quad (122)$$

where (a) follows from  $\|\mathbf{x} + \mathbf{y}\|_\infty^2 \leq 2\|\mathbf{x}\|_\infty^2 + 2\|\mathbf{y}\|_\infty^2$ . ■

Now, let  $\mathbf{p}^*$  be the optimal solution of (P4). Notice that,

$$\begin{aligned} -\frac{1}{2\alpha}\|\nabla f_t(\mathbf{p}(t-1))\|_\infty^2 &\leq_{(a)} \nabla f_t(\mathbf{p}(t-1))^\top (\mathbf{p}(t) - \mathbf{p}(t-1)) + \alpha D(\mathbf{p}(t) \| \mathbf{p}(t-1)) \\ &\leq_{(b)} \nabla f_t(\mathbf{p}(t-1))^\top (\mathbf{p}^* - \mathbf{p}(t-1)) + \alpha D(\mathbf{p}^* \| \mathbf{p}(t-1)) - \alpha D(\mathbf{p}^* \| \mathbf{p}(t)) \\ &\leq_{(c)} f_t(\mathbf{p}^*) - f_t(\mathbf{p}(t-1)) + \alpha D(\mathbf{p}^* \| \mathbf{p}(t-1)) - \alpha D(\mathbf{p}^* \| \mathbf{p}(t)), \end{aligned} \quad (123)$$

where (a) follows from Lemma 15, (b) follows from Lemma 5 with  $\mathcal{G} = [0, \infty)^n$ ,  $\mathcal{C} = \mathcal{I}$ , and  $\omega(\mathbf{x}) = \sum_{j=1}^n x_j \ln(x_j)$  (Notice that from (116),  $\mathbf{p}(t) \in \mathcal{I}^0$  which is required for the lemma), and (c) follows from sub-gradient inequality for the convex function  $f_t$ . Summing the above inequality for  $t$  from 1 to  $T$ , and dividing by  $T$ , we have,

$$\begin{aligned} -\sum_{t=1}^T \frac{\|\nabla f_t(\mathbf{p}(t-1))\|_\infty^2}{2\alpha T} &\leq \frac{\sum_{t=1}^T f_t(\mathbf{p}^*)}{T} - \frac{\sum_{t=1}^T f_t(\mathbf{p}(t-1))}{T} + \frac{\alpha D(\mathbf{p}^* \| \mathbf{p}(0))}{T} - \frac{\alpha D(\mathbf{p}^* \| \mathbf{p}(T))}{T} \\ &\leq \frac{\sum_{t=1}^T f_t(\mathbf{p}^*)}{T} - \frac{\sum_{t=1}^T f_t(\mathbf{p}(t-1))}{T} + \frac{\alpha}{T} \left( \sum_{k=1}^n p_k^* \ln\left(\frac{p_k(T)}{p_k(0)}\right) \right) \\ &\leq_{(a)} \frac{\sum_{t=1}^T f_t(\mathbf{p}^*)}{T} - \frac{\sum_{t=1}^T f_t(\mathbf{p}(t-1))}{T} + \frac{\alpha}{T} \max\left(\ln\left(\frac{1}{p_k(0)}\right); 1 \leq k \leq n\right) \\ &\leq_{(b)} \frac{\sum_{t=1}^T f_t(\mathbf{p}^*)}{T} - \frac{\sum_{t=1}^T f_t(\mathbf{p}(t-1))}{T} + \frac{\alpha}{T} \ln(n), \end{aligned} \quad (124)$$

where (a) follows from the facts that natural logarithm is an increasing function in  $(0, \infty)$ , and  $p_k(T) \leq 1$  for all  $k$ , and (b) follows from  $p_k(0) = 1/n$ . After taking expectations of both sides and some rearrangements, we have

$$-\sum_{t=1}^T \frac{\mathbb{E}\{\|\nabla f_t(\mathbf{p}(t-1))\|_\infty^2 | \mathbf{Z}\}}{2\alpha T} - \frac{\alpha}{T} \ln(n)$$

$$\leq \mathbb{E} \left\{ \frac{\sum_{t=1}^T f_t(\mathbf{p}^*)}{T} \middle| \mathbf{Z} \right\} - \mathbb{E} \left\{ \frac{\sum_{t=1}^T f_t(\mathbf{p}(t-1))}{T} \middle| \mathbf{Z} \right\}. \quad (125)$$

Substituting the results from Lemma 16 in (125), we have that,

$$-\frac{C}{2\alpha} - \frac{\alpha}{T} \ln(n) \leq \mathbb{E} \left\{ \frac{\sum_{t=1}^T f_t(\mathbf{p}^*)}{T} \middle| \mathbf{Z} \right\} - \mathbb{E} \left\{ \frac{\sum_{t=1}^T f_t(\mathbf{p}(t-1))}{T} \middle| \mathbf{Z} \right\}, \quad (126)$$

where

$$C = 2\|\mathbf{E}\|_\infty^2 + \frac{1}{2}\mathbb{E} \{ \|\boldsymbol{\Omega}\|_\infty^2 | \mathbf{Z} \}. \quad (127)$$

But notice that,

$$\mathbb{E} \left\{ \frac{\sum_{t=1}^T f_t(\mathbf{p}^*)}{T} \right\} = -\frac{1}{T} \sum_{t=1}^T \mathbb{E} \{ f_t(\mathbf{p}^*) | \mathbf{Z} \} = \frac{1}{T} \sum_{t=1}^T f(\mathbf{p}^*) = -f^*. \quad (128)$$

Also, we have that,

$$\begin{aligned} \mathbb{E} \left\{ \frac{\sum_{t=1}^T f_t(\mathbf{p}(t-1))}{T} \middle| \mathbf{Z} \right\} &= \frac{\sum_{t=1}^T \mathbb{E} \{ f_t(\mathbf{p}(t-1)) | \mathbf{Z} \}}{T} \\ &= \frac{\sum_{t=1}^T \mathbb{E} \{ \mathbb{E}_{\boldsymbol{\Omega}(t)} \{ f_t(\mathbf{p}(t-1)) | \boldsymbol{\Omega}(\tau) \text{ for } t \neq \tau \} | \mathbf{Z} \}}{T} \\ &=_{(a)} \frac{\sum_{t=1}^T \mathbb{E} \{ -f(\mathbf{p}(t-1)) | \mathbf{Z} \}}{T} \\ &= -\mathbb{E} \left\{ \frac{\sum_{t=1}^T f(\mathbf{p}(t-1)) | \mathbf{Z} \right\} \geq -\mathbb{E} \left\{ f \left( \frac{\sum_{t=1}^T \mathbf{p}(t-1)}{T} \right) \middle| \mathbf{Z} \right\}, \end{aligned} \quad (129)$$

where (a) follows since  $\boldsymbol{\Omega}(t)$  is independent of  $\boldsymbol{\Omega}(\tau)$  for  $\tau \neq t$  and  $\mathbf{p}(t-1)$  is a function of  $\boldsymbol{\Omega}(\tau)$  for  $\tau < t$  and hence is independent of  $\boldsymbol{\Omega}(t)$ . The last inequality follows from Jensen's inequality since  $f$  is a concave function. Substituting, in (126), we have the desired result. ■

### B. Case $a = 1$

When  $a = 1$ , the set  $\mathcal{G}^A$  is more complex than the simplex  $\mathcal{I}$ . The next two lemmas compute for each  $\mathbf{p} \in \mathcal{I}$ , the largest  $q \in \mathbb{R}$ , such that  $(q, \mathbf{p}) \in \mathcal{G}^A$ . For simplicity of exposition, this subsection assumes  $W_1$  has a continuous CDF  $F_{W_1} : \mathbb{R} \rightarrow [0, 1]$ .

*Lemma 17:* If  $(q, \mathbf{p}) \in \mathcal{G}^A$ , then we have,

$$q \leq \mathbb{E} \{ W_1 | F_{W_1}(W_1) \geq 1 - p_1 \} p_1. \quad (130)$$

*Proof:* See Appendix I ■



*Lemma 18:* Consider any  $\mathbf{p} \in \mathcal{I}$ . Let  $q = \mathbb{E}\{W_1 | F_{W_1}(W_1) \geq 1 - p_1\} p_1$ . Then  $(q, \mathbf{p}) \in \mathcal{G}^A$ . In particular, there exists  $g^A \in \mathcal{S}^A$  such that,

$$\begin{aligned} p_k &= \mathbb{E}\{\mathbb{1}_{g^A(U^A, W_1, \mathbf{Z})=k} | \mathbf{Z}\} \text{ for } 1 \leq k \leq n, \\ q &= \mathbb{E}\{W_k \mathbb{1}_{g^A(U^A, W_1, \mathbf{Z})=1} | \mathbf{Z}\}. \end{aligned} \quad (131)$$

*Proof:* See Appendix J ■

Given that  $W_1$  has a continuous CDF, for a fixed  $\mathbf{p}$ , it is optimal to choose the strategy given by Lemma 18 due to Theorem 2-2. Hence we can simply focus on solving the problem,

$$\begin{aligned} \text{(P5): maximize } & f(\mathbb{E}\{W_1 | F_{W_1}(W_1) \geq 1 - p_1\} p_1, \mathbf{p}_{2:n}) \\ \text{subject to } & \mathbf{p} \in \mathcal{I}, \end{aligned} \quad (132)$$

It can be shown that  $\mathbb{E}\{W_1 | F_{W_1}(W_1) \geq 1 - p_1\} p_1$  is concave and non-increasing in  $p_1$ . Hence (P5) is also a convex optimization problem since  $f$  is entry-wise, non-decreasing, and concave. Hence, when  $\mathbb{E}\{W_1 | F_{W_1}(W_1) \geq 1 - p_1\} p_1$  has a bounded subgradient at  $p_1 = 0$ , its subgradients are bounded for all  $p_1 \in [0, 1]$ , and (P5) can be solved using the mirror-descent algorithm described for  $a = 0$ . When subgradients at  $p_1 = 0$  are not bounded, extra care should be taken, such as by restricting the search to a subregion of  $\mathcal{I}$ .

## APPENDIX I

### PROOF OF LEMMA 17

Let  $\tau_1 \in \mathbb{R}$ , be such that,  $F_{W_1}(W_1) = 1 - p_1$ . Since the CDF of  $W_1$  is assumed to be continuous, the intermediate value theorem guarantees the existence of such a  $\tau_1$ . Let  $f_1 : \mathbb{R} \rightarrow \mathbb{R}$  be defined such that,

$$f_1(x) = \begin{cases} 0 & \text{if } x \leq \tau_1, \\ 1 & \text{otherwise.} \end{cases} \quad (133a)$$

$$(133b)$$

Hence, we have that  $p_1 \mathbb{E}\{W_1 | F_{W_1}(W_1) \geq 1 - p_1\} = \mathbb{E}\{W_1 f_1(W_1)\} = \mathbb{E}\{W_1 f_1(W_1) | \mathbf{Z}\}$ , where the last equality follows from the fact that  $\mathbf{Z}$  and  $W_1$  are independent.

Recall the definition of strategy function in (73). Let  $P^A$  be the strategy function of the strategy relating to  $(q, \mathbf{p})$  (See (73)). Hence,

$$\begin{aligned} q - p_1 \mathbb{E}\{W_1 | F_{W_1}(W_1) \geq 1 - p_1\} &= \mathbb{E}\{W_1 (P^A(W_1, \mathbf{Z}) - f_1(W_1)) | \mathbf{Z}\} \\ &=_{(a)} \Pr\{W_1 \leq \tau_1\} \mathbb{E}\{W_1 (P^A(W_1, \mathbf{Z}) - f_1(W_1)) | \mathbf{Z}, W_1 \leq \tau_1\} \end{aligned}$$

$$\begin{aligned}
& + \Pr\{W_1 > \tau_1\} \mathbb{E}\{W_1(P^A(W_1, \mathbf{Z}) - f_1(W_1)) | \mathbf{Z}, W_1 > \tau_1\} \\
& = \Pr\{W_1 \leq \tau_1\} \mathbb{E}\{W_1 P^A(W_1, \mathbf{Z}) | \mathbf{Z}, W_1 \leq \tau_1\} \\
& + \Pr\{W_1 > \tau_1\} \mathbb{E}\{W_1(P^A(W_1, \mathbf{Z}) - 1) | \mathbf{Z}, W_1 > \tau_1\} \\
& \leq \Pr\{W_1 \leq \tau_1\} \mathbb{E}\{\tau_1 P^A(W_1, \mathbf{Z}) | \mathbf{Z}, W_1 \leq \tau_1\} \\
& + \Pr\{W_1 > \tau_1\} \mathbb{E}\{\tau_1(P^A(W_1, \mathbf{Z}) - 1) | \mathbf{Z}, W_1 > \tau_1\} \\
& = \Pr\{W_1 \leq \tau_1\} \mathbb{E}\{\tau_1 P^A(W_1, \mathbf{Z}) | \mathbf{Z}, W_1 \leq \tau_1\} \\
& + \Pr\{W_1 > \tau_1\} \mathbb{E}\{\tau_1 P^A(W_1, \mathbf{Z}) | \mathbf{Z}, W_1 > \tau_1\} - \Pr\{W_1 > \tau_1\} \mathbb{E}\{\tau_1 | \mathbf{Z}, W_1 \leq \tau_1\} \\
& =_{(b)} \mathbb{E}\{\tau_1 P^A(W_1, \mathbf{Z}) | \mathbf{Z}\} - \Pr\{W_1 > \tau_1\} \mathbb{E}\{\tau_1 | \mathbf{Z}, W_1 \leq \tau_1\} = \tau_1 p_1 - p_1 \tau_1 = 0,
\end{aligned}$$

where (a) and (b) follow from the total probability law. Hence,  $q \leq p_1 \mathbb{E}\{W_1 | F_{W_1}(W_1) \geq 1 - p_1\}$ .

## APPENDIX J

### PROOF OF LEMMA 18

Similar to Appendix I, define  $\tau_1 \in \mathbb{R}$ , be such that,  $F_{W_1}(W_1) = 1 - p_1$ . Now consider the strategy for A of choosing 1 whenever  $\tau_1 \leq W_1$  and  $\{2, \dots, n\}$  with probabilities  $\{p_2, \dots, p_n\}$  otherwise. Note that  $\Pr\{\alpha^A = 1 | \mathbf{Z}\} = \Pr\{\tau_1 \leq W_1\} = p_1$ . Moreover, notice that

$$\mathbb{E}\{W_1 \mathbb{1}_{\alpha^A=1} | \mathbf{Z}\} = p_1 \mathbb{E}\{W_1 | \alpha^A = 1, \mathbf{Z}\} = p_1 \mathbb{E}\{W_1 | F_{W_1}(W_1) \geq 1 - p_1\} = q, \quad (134)$$

which follows since  $\mathbf{Z}$  and  $W_1$  are independent. Hence,  $(q, \mathbf{p}) \in \mathcal{G}^A$ , as desired.

## REFERENCES

- [1] K. Akkarajitsakul, E. Hossain, D. Niyato, and D. I. Kim, "Game theoretic approaches for multiple access in wireless networks: A survey," *IEEE Communications Surveys & Tutorials*, vol. 13, no. 3, pp. 372–395, 2011.
- [2] E. Aryafar, A. Keshavarz-Haddad, M. Wang, and M. Chiang, "Rat selection games in hetnets," in *2013 Proceedings IEEE INFOCOM*, 2013, pp. 998–1006.
- [3] M. Felegyhazi, M. Cagalj, S. S. Bidokhti, and J.-P. Hubaux, "Non-cooperative multi-radio channel allocation in wireless networks," in *IEEE INFOCOM 2007 - 26th IEEE International Conference on Computer Communications*, 2007, pp. 1442–1450.
- [4] B. Li, Q. Qu, Z. Yan, and M. Yang, "Survey on ofdma based mac protocols for the next generation wlan," in *2015 IEEE Wireless Communications and Networking Conference Workshops (WCNCW)*, 2015, pp. 131–135.
- [5] R. W. Rosenthal, "A class of games possessing pure-strategy Nash equilibria," *International Journal of Game Theory*, vol. 2, pp. 65–67, 1973.

- [6] E. Nikolova and N. E. Stier-Moses, “Stochastic selfish routing,” in *Algorithmic Game Theory*, G. Persiano, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 314–325.
- [7] H. Angelidakis, D. Fotakis, and T. Lianas, “Stochastic congestion games with risk-averse players,” in *Lecture Notes in Computer Science*, 10 2013.
- [8] C. Zhou, T. H. Nguyen, and H. Xu, “Algorithmic information design in multi-player games: Possibilities and limits in singleton congestion,” in *Proceedings of the 23rd ACM Conference on Economics and Computation*, ser. EC '22. New York, NY, USA: Association for Computing Machinery, 2022, p. 869. [Online]. Available: <https://doi.org/10.1145/3490486.3538238>
- [9] M. Castiglioni, A. Celli, A. Marchesi, and N. Gatti, “Signaling in Bayesian network congestion games: The subtle power of symmetry,” in *AAAI Conference on Artificial Intelligence*, 2020.
- [10] M. Wu, J. Liu, and S. Amin, “Informational aspects in a class of Bayesian congestion games,” in *2017 American Control Conference (ACC)*, 2017, pp. 3650–3657.
- [11] V. Syrgkanis, “The complexity of equilibria in cost sharing games,” in *Internet and Network Economics: 6th International Workshop, WINE 2010, Stanford, CA, USA, December 13-17, 2010. Proceedings 6*. Springer, 2010, pp. 366–377.
- [12] M. Zinkevich, “Online convex programming and generalized infinitesimal gradient ascent,” in *Proceedings of the Twentieth International Conference on International Conference on Machine Learning*, ser. ICML'03. AAAI Press, 2003, p. 928–935.
- [13] H. Yu, M. Neely, and X. Wei, “Online convex optimization with stochastic constraints,” in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017.
- [14] M. J. Neely, *Stochastic Network Optimization with Application to Communication and Queueing Systems*. Morgan & Claypool, 2010.
- [15] D. Monderer and L. S. Shapley, “Potential games,” *Games and Economic Behavior*, vol. 14, no. 1, pp. 124–143, 1996. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0899825696900445>
- [16] S. Chien and A. Sinclair, “Convergence to approximate Nash equilibria in congestion games,” *Games and Economic Behavior*, vol. 71, no. 2, pp. 315–327, 2011. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0899825609001110>
- [17] K. Bhawalkar, M. Gairing, and T. Roughgarden, “Weighted congestion games: Price of anarchy, universal worst-case examples, and tightness,” in *Lecture Notes in Computer Science*, 09 2010, pp. 17–28.
- [18] I. Milchtaich, “Congestion games with player-specific payoff functions,” *Games and Economic Behavior*, vol. 13, no. 1, pp. 111–124, 1996. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0899825696900275>
- [19] H. Ackermann, P. W. Goldberg, V. S. Mirrokni, H. Röglin, and B. Vöcking, “A Unified Approach to Congestion Games and Two-Sided Markets,” *Internet Mathematics*, vol. 5, no. 4, pp. 439 – 458, 2008. [Online]. Available: <https://doi.org/>
- [20] D. Fotakis, S. Kontogiannis, E. Koutsoupias, M. Mavronicolas, and P. Spirakis, “The structure and complexity of Nash equilibria for a selfish routing game,” *Theoretical Computer Science*, vol. 410, no. 36, pp. 3305–3326, 2009, graphs, Games and Computation: Dedicated to Professor Burkhard Monien on the Occasion of his 65th Birthday. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0304397508000261>
- [21] M. Gairing, T. Lücking, M. Mavronicolas, and B. Monien, “Computing Nash equilibria for scheduling on restricted parallel links,” in *Proceedings of the Thirty-Sixth Annual ACM Symposium on Theory of Computing*, ser. STOC '04. New York, NY, USA: Association for Computing Machinery, 2004, p. 613–622. [Online]. Available: <https://doi.org/10.1145/1007352.1007446>

- [22] D. Acemoglu, A. Makhdoumi, A. Malekian, and A. Ozdaglar, “Informational braess’ paradox: The effect of information on traffic congestion,” *Operations Research*, vol. 66, no. 4, pp. 893–917, 2018. [Online]. Available: <https://doi.org/10.1287/opre.2017.1712>
- [23] S. Le, Y. Wu, and M. Toyoda, “A congestion game framework for service chain composition in NFV with function benefit,” *Inf. Sci.*, vol. 514, no. C, p. 512–522, apr 2020. [Online]. Available: <https://doi.org/10.1016/j.ins.2019.11.015>
- [24] L. Zhang, K. Gong, and M. Xu, “Congestion control in charging stations allocation with Q-learning,” *Sustainability*, vol. 11, no. 14, 2019. [Online]. Available: <https://www.mdpi.com/2071-1050/11/14/3900>
- [25] E. Anshelevich, A. Dasgupta, J. Kleinberg, E. Tardos, T. Wexler, and T. Roughgarden, “The price of stability for network design with fair cost allocation,” in *45th Annual IEEE Symposium on Foundations of Computer Science*, 2004, pp. 295–304.
- [26] I. Caragiannis, M. Flammini, C. Kaklamani, P. Kanellopoulos, and L. Moscardelli, “Tight bounds for selfish and greedy load balancing,” *Algorithmica*, vol. 58, pp. 311–322, 01 2006.
- [27] F. Zhang and M. M. Wang, “Stochastic congestion game for load balancing in mobile-edge computing,” *IEEE Internet of Things Journal*, vol. 8, no. 2, pp. 778–790, 2021.
- [28] M. Liu, S. H. A. Ahmad, and Y. Wu, “Congestion games with resource reuse and applications in spectrum sharing,” in *2009 International Conference on Game Theory for Networks*, 2009, pp. 171–179.
- [29] M. Liu and Y. Wu, “Spectrum sharing as congestion games,” in *2008 46th Annual Allerton Conference on Communication, Control, and Computing*. IEEE, 2008, pp. 1146–1153.
- [30] M. Ibrahim, K. Khawam, and S. Tohme, “Congestion games for distributed radio access selection in broadband networks,” in *2010 IEEE Global Telecommunications Conference GLOBECOM 2010*, 2010, pp. 1–5.
- [31] J.-B. Seo and H. Jin, “Two-user noma uplink random access games,” *IEEE Communications Letters*, vol. 22, no. 11, pp. 2246–2249, 2018.
- [32] —, “Revisiting two-user s-aloha games,” *IEEE Communications Letters*, vol. 22, no. 6, pp. 1172–1175, 2018.
- [33] I. Malanchini, M. Cesana, and N. Gatti, “Network selection and resource allocation games for wireless access networks,” *IEEE Transactions on Mobile Computing*, vol. 12, no. 12, pp. 2427–2440, 2013.
- [34] R. Trestian, O. Ormond, and G.-M. Muntean, “Game theory-based network selection: Solutions and challenges,” *IEEE Communications Surveys & Tutorials*, vol. 14, no. 4, pp. 1212–1231, 2012.
- [35] T. Quint and M. Shubik, “A model of migration,” *Working Paper*, 1994. [Online]. Available: <https://elischolar.library.yale.edu/cowles-discussion-paper-series/1331>
- [36] N. Nisan, T. Roughgarden, E. Tardos, and V. V. Vazirani, *Algorithmic Game Theory*. Cambridge University Press, 2007.
- [37] X. Wei, H. Yu, and M. J. Neely, “Online primal-dual mirror descent under stochastic constraints,” *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 4, no. 2, jun 2020. [Online]. Available: <https://doi.org/10.1145/3392157>
- [38] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, “Robust stochastic approximation approach to stochastic programming,” *Society for Industrial and Applied Mathematics*, vol. 19, pp. 1574–1609, 01 2009.