

Probably Correct Optimal Stable Matching for Two-Sided Markets Under Uncertainty

Andreas Athanasopoulos

University of Neuchâtel

Neuchâtel, Switzerland

andreas.athanasopoulos@unine.ch

Anne-Marie George

University of Oslo

Oslo, Norway

annemage@ifi.uio.no

Christos Dimitrakakis

University of Neuchâtel

Neuchâtel, Switzerland

christos.dimitrakakis@gmail.com

ABSTRACT

We consider a learning problem for the stable marriage model under unknown preferences for the left side of the market. We focus on the centralized case, where at each time step, an online platform matches the agents, and obtains a noisy evaluation reflecting their preferences. Our aim is to quickly identify the stable matching that is left-side optimal, rendering this a pure exploration problem with bandit feedback. We specifically aim to find Probably Correct Optimal Stable Matchings and present several bandit algorithms to do so. Our findings provide a foundational understanding of how to efficiently gather and utilize preference information to identify the optimal stable matching in two-sided markets under uncertainty. An experimental analysis on synthetic data complements theoretical results on sample complexities for the proposed methods.

KEYWORDS

Two-sided matching markets; Stable matching; Pure exploration

ACM Reference Format:

Andreas Athanasopoulos, Anne-Marie George, and Christos Dimitrakakis. 2025. Probably Correct Optimal Stable Matching for Two-Sided Markets Under Uncertainty. In *Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025)*, Detroit, Michigan, USA, May 19 – 23, 2025, IFAAMAS, 9 pages.

1 INTRODUCTION

We explore the problem of two-sided matching markets, where two distinct groups of agents must be matched with one another [21]. Such markets have diverse real-world applications, ranging from labor markets, like online crowd-sourcing platforms such as Amazon Mechanical Turk, to online dating services [9, 10, 20]. In these contexts, the challenge lies in designing a matching algorithm that respects the preferences of both sets of agents.

In their seminal work, Gale and Shapley introduce the concept of stable matchings [9]. Here, agents are matched one-to-one according to ordinal preferences such that no two agents have the incentive to deviate from the proposed solution. Gale and Shapley’s Deferred Acceptance algorithm guarantees a stable matching that is optimal for one side of the market. This optimal stable matching always exists and is unique, based on one side’s preferences. Many works have considered solutions for this and related problems of

finding matchings under known preferences [18, 23]. However, in online settings, agents often lack certainty about their preferences.

More recently, research has focused on cases where agents’ preferences are learned over time. This line of work, initiated by Das and Kamenica, frames the problem as an online decision-making challenge, likening it to a *multi-agent multi-armed bandit* problem where agents must compete for potential matches. In the centralized version of this problem, a platform matches agents according to their preferences, and the agents receive noisy feedback that updates those preferences. In this setting, there is an inherent tension between exploration and exploitation: the platform must decide whether to exploit current estimates of the best stable match or continue exploring to refine agents’ preferences. Liu et al. formalize this dilemma using regret measures that balance this trade-off.

In this work, we focus on the centralized version of the stable marriage model, but approach the problem from a pure exploration standpoint. Our primary goal is to *efficiently* learn the optimal stable matching w.r.t. the side of the market with uncertain preferences. This perspective emerges from the observation that algorithms designed for the regret-minimization setting, which explore potential stable matches, can hinder the exploration of critical agent pairs that contribute to stability. Accordingly, we move beyond the regret-minimization framework to address the problem:

How can we efficiently identify the optimal stable matching with high probability?

Our contributions can be summarised as follows:

- (1) **Action Space.** We do not limit our algorithms to stable matchings during exploration, so as to be able to learn player preferences. This is because constraining ourselves to stable matchings may never lead to discovering the optimal stable match (Section 3.2).
- (2) **Solution concept.** We introduce the concept of a *probably correct optimal stable matching* (Section 3.3). This special case of probably approximately correct (PAC) solutions requires the output matching to be optimal with high probability.
- (3) **Uniform exploration strategy.** We begin our analysis with an algorithm that uniformly samples all available agent pairs, similar to the Explore-Then-Commit (ETC) strategy. We demonstrate that this can produce the optimal stable matching with high probability and provide a bound on its sample complexity (Section 4).
- (4) **Action elimination algorithms.** Next, we explore an algorithm based on the concept of action elimination (Section 5), which improves sample efficiency and reduces dependence on instance-specific parameters. Additionally, we enhance sample complexity by modifying the stopping criterion, enabling the algorithm to terminate when sufficient information is gathered (Section 6).
- (5) **Adaptive sampling.** We study a strategy that adaptively samples agent pairs, improving the complexity in practice (Section 7).



This work is licensed under a Creative Commons Attribution International 4.0 License.

Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025), Y. Vorobeychik, S. Das, A. Nowé (eds.), May 19 – 23, 2025, Detroit, Michigan, USA. © 2025 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org).

(6) Experimental evaluation. Finally, we compare the performance of our algorithms on several synthetic instances (Section 8).

2 RELATED WORK

The *stochastic* MAB problem, introduced by [24], has been extensively studied in the past few decades, becoming a fundamental tool for analyzing online decision-making problems under uncertainty. In its classic version, a learner selects one arm in each round and receives a numeric reward drawn from an unknown distribution associated with the selected arm. The goal of maximising the learner’s cumulative reward requires a balanced solution of the exploration-exploitation dilemma: finding the best arm requires exploration, while receiving good rewards requires exploitation. Another popular setting is the pure exploration problem [6, 7, 19], where the learner’s objective is to identify the best action (arm to pull) regardless of the cost incurred by choosing suboptimal actions but with a minimal number of samples. This problem is usually studied in the Probably Approximately Correct (PAC) setting, where the learner aims to identify the approximately best arm with high probability. This paper focuses on the pure exploitation setting and the presented solutions like ETC and UCB strategies. In contrast to the PAC setting, we aim at identifying the best action with high probability, not just an approximation. For a comprehensive overview of MAB, we refer to the book by Lattimore and Szepesvári.

The stable marriage problem, is a well studied problem in two-sided markets in which agents from the two sides are matched such that no pair of agents prefers to deviate from the matching [9, 18], making the solution a *stable matching*. In earlier work on stable matchings, preferences of agents in the market are assumed to be known. The problem of two-sided markets with unknown preferences has recently attracted considerable attention in various settings. Das and Kamenica were the first to introduce the marriage problem with unknown preferences on both sides of the market, framing it as a MAB problem. In their work [5], they empirically studied several algorithms in specific preference settings. A follow-up work [16] studied a variant of the problem where one side of the market has unknown preferences and first introduced the notion of Player Stable Optimal/Pessimal regret for the agents in the market. They studied an ETC-like algorithm, proving sublinear bounds for both notions of regret, but requiring additional knowledge of the reward differences of the arms. They also study a UCB-like approach with sublinear bounds on the Player Pessimal Stable regret. However, they identified fundamental issues with their UCB algorithm in achieving sublinear optimal stable regret for the agents. Another research direction studied the decentralized setting where the agents act independently in the market [3, 14, 17, 22]. Kong et al. introduced a Thompson sampling algorithm for the decentralized market, and they also highlight issues with their algorithm in achieving sublinear optimal/pessimal stable regret in the centralized setting. Overall, developing algorithms with sublinear optimal regret remains an open problem. To bridge this gap, other lines of research study the centralized problem with transferable utilities to guarantee stability [4, 12]. In this work, we instead focus on approaching the problem in a pure exploration setting, aiming to efficiently identify a correct optimal stable matching with high

probability. To the best of our knowledge, the pure exploration problem has not been explored in this setting.

3 PROBLEM SETTING

In this section, we first introduce all necessary definitions of agents, their preferences and stability of matchings for the classic two-sided matching markets. We then describe the problem of learning an optimal stable matching in two-sided matching markets with unknown preferences.

3.1 Stable Matchings in Two-Sided Markets

Agents: We consider two distinct sets of agents, *players* $\mathcal{P} = \{p_1, \dots, p_N\}$ and *arms* $\mathcal{A} = \{a_1, \dots, a_K\}$, each having N and K elements, respectively, where $N \leq K$. Let the set of all agents be $\mathfrak{A} = \mathcal{P} \cup \mathcal{A}$. We will also refer to the players as the *left* and the arms as the *right* side of the market.

Preferences: Each agent’s preferences are represented as a complete list of agents from the opposite side. More specifically, the preferences π_p of a player p are given as a permutation over the arms, i.e., $\pi_p = (a_{p_1}, \dots, a_{p_K}) \in \mathbf{P}(\mathcal{A})$, where $\mathbf{P}(\mathcal{A})$ is the set of permutations over \mathcal{A} . We say player p prefers arm $\pi_p(i)$ to arm $\pi_p(j)$ if $i < j$. Sometimes, we express these preference by the *order relation* \succ_{π_p} where $\pi_p(i) \succ_{\pi_p} \pi_p(j)$ if $i < j$. Here we might omit the subscript, simply writing \succ , if it is clear from the context. Further, we define the *rank of an arm* a w.r.t. π_p as the position of a in p ’s preference list, i.e., $r_p(a) := |\{a_j \in \mathcal{A} : a_j \succ_{\pi_p} a\}|$. Similarly, we can define preferences π_a , the corresponding order relation \succ_{π_a} , and rank function r_a of an arm $a \in \mathcal{A}$ over the set of players \mathcal{P} .

Matchings: In the stable marriage model with players \mathcal{P} and arms \mathcal{A} , a *matching* is a set of player-arm pairs $m \subseteq \mathcal{P} \times \mathcal{A}$ that are pairwise disjoint, representing the pairs of agents that are *matched*. We say an agent $i \in \mathfrak{A}$ is *unmatched* under m if there exists no pair in m that involves i . By a slight abuse of notation, we define the equivalent functional representation $m : \mathfrak{A} \rightarrow \mathfrak{A} \cup \{\perp\}$ of a matching m , where $m(i) = j$ and $m(j) = i$ for the pair $(i, j) \in m$, and $m(i) = \perp$ if $i \in \mathfrak{A}$ is unmatched. Let \mathcal{M} denote the set of all possible matchings.

Stability: In order for a matching m to align with the agents’ preferences, Gale and Shapley proposed stability [9] as a notion of equilibrium in the market. A matching m is *stable* if there is no pair of agents who prefer to be matched with each other compared to their match under m . More formally, m is a stable matching, if there exists no *blocking pair* $(p, a) \in \mathcal{P} \times \mathcal{A}$, i.e., no pair $(p, a) \notin m$ such that (1) $a \succ_{\pi_p} m(p)$ or p is unmatched, and (2) $p \succ_{\pi_a} m(a)$ or a is unmatched.

In Gale and Shapley’s Deferred Acceptance (DA) algorithm [9], the agents on one side of the market sequentially propose to the other side, while the other side is temporarily matched with the most preferred agents so far until all agents are matched. They demonstrate that not only does a stable matching always exist for any instance of the marriage problem, but also that multiple stable matchings can exist; we denote this *collection of stable matchings* by \mathcal{S} . The stable matching m_s^* produced by the algorithm is optimal for the proposing side of the market and pessimal for the side that received the proposals, in the sense that the agents of the sets are matched with their most/least preferred partner among any feasible

stable matching in \mathcal{S} . In the remainder of this paper, we refer to the unique optimal stable matching, m_s^* , w.r.t. the players' preferences as the *optimal stable matching*.

3.2 Two-Sided Markets with Unknown Player Preferences

In this work, we consider a setting where the players $p \in \mathcal{P}$ are initially uncertain about their preferences π_p . Then the players gradually learn their preferences through noisy feedback from matchings that are imposed by a centralised matching algorithm A.

More formally, the learning process is performed in T rounds where in each round t the algorithm A selects a matching m_t for the agents. Then, each player $p \in \mathcal{P}$ receives a stochastic reward $X_{p,m_t(p)}^t \in [0, 1]$, distributed according to an unknown distribution $\mathbb{P}_{p,m_t(p)}$ with mean $\mu_{p,m_t(p)}$. We assume that a player p truly prefers arm a_1 to arm a_2 if $\mu_{p,a_1} > \mu_{p,a_2}$, while we denote the difference of their expected rewards by $\Delta_{p,a_1,a_2} = |\mu_{p,a_1} - \mu_{p,a_2}| > 0$.

After each round t , the players update their preferences $\hat{\pi}_p(t)$ based on the estimate of the expected rewards $\hat{\mu}_{p,a}(t)$ (sample mean) and report it to the algorithm. The objective of the algorithm is to efficiently identify the optimal stable matching m_s^* with respect to the (initially unknown) preferences π_p with high probability. Thus, at each round t , the algorithm aims to select a matching $m_t \in \mathcal{M}$ to gather relevant information.

To summarize the process, at every round $t \in \{1, \dots, T\}$:

- (1) The algorithm A selects a matching $m_t \in \mathcal{M}$.
- (2) The players $p \in \mathcal{P}$ receive a reward $X_{p,m_t(p)}^t$.
- (3) The players update $\hat{\mu}_{p,m_t(p)}(t)$ and report it to A.
- (4) The algorithm A might terminate the process (according to a stopping criterion) and return a matching.

The returned matching should be the true optimal stable matching m_s^* with high probability.

Note that here we do not require the action selected by A to be a *stable* matching. Since preferences of agents are uncertain, we cannot be certain which matchings are stable. Yet related work restricts A to select only matchings that are stable with respect to confidence bounds [4, 5, 16]. Because we can compute confidence bounds on the estimated means μ_p , we could identify blocking pairs with high probability. However, the following example, which Liu et al. use to demonstrate that their UCB-like algorithm fails to achieve sublinear player-optimal stable regret [16], shows that excluding actions of matchings with high-probability blocking pairs will in some cases not allow us to identify the optimal stable matching. This motivates our setting in which the action space of algorithm A is the set of all matchings \mathcal{M} disregarding stability.

Example 1 (Liu et al. [16]). Consider a market with three agents on each side and the following true agents' preferences.

| Players | Arms |
|---------------------------------------|---------------------------------------|
| $\pi_{p_1} : a_1 \succ a_2 \succ a_3$ | $\pi_{a_1} : p_2 \succ p_3 \succ p_1$ |
| $\pi_{p_2} : a_2 \succ a_1 \succ a_3$ | $\pi_{a_2} : p_1 \succ p_2 \succ p_3$ |
| $\pi_{p_3} : a_3 \succ a_1 \succ a_2$ | $\pi_{a_3} : p_3 \succ p_1 \succ p_2$ |

Here, the only stable matchings are the player-optimal stable matching $m_{s1} = \{(p_1, a_1), (p_2, a_2), (p_3, a_3)\}$ and the arm-optimal stable matching $m_{s2} = \{(p_1, a_2), (p_2, a_1), (p_3, a_3)\}$ which is player-pessimal.

Now assume that players p_1 and p_2 are certain about their preferences and correctly report them, while player p_3 is uncertain about the order of a_1 and a_3 , but p_3 's current estimates of the sample means $\hat{\mu}_{p_3,a_1}$, $\hat{\mu}_{p_3,a_2}$ and $\hat{\mu}_{p_3,a_3}$ yield $\hat{\pi}_{p_3} : a_1 \succ a_3 \succ a_2$. Under these preferences, the only stable matching is m_{s2} . Further, the only matchings that could be stable w.r.t any possible preference $\hat{\pi}_{p_3} \in \mathcal{P}(\mathcal{A})$ of p_3 , are m_{s1} or m_{s2} .

In this situation, if the matching algorithm A is only permitted to select potentially stable matchings, i.e., m_{s1} or m_{s2} , the player p_3 will receive a sample from \mathbb{P}_{p_3,a_3} . This serves p_3 to grow more certain about the estimate $\hat{\mu}_{p_3,a_3}$. However, the uncertainty about $\hat{\mu}_{p_3,a_1}$ persists. Thus, if $\hat{\mu}_{p_3,a_1} > \mu_{p_3,a_3}$ then with high probability the updated estimated preference order $\hat{\pi}_{p_3}$ will remain $a_1 \succ a_3 \succ a_2$, i.e., yielding the same situation as before. Thus, no matter the number of samples, with high probability the returned matching is m_{s2} — the player-pessimal stable matching!

In contrast, if A is permitted to sample any matching, it could also sample a matching that matches p_3 to a_1 in order to grow more certain about the estimate $\hat{\mu}_{p_3,a_1}$ and eventually determining $\hat{\mu}_{p_3,a_3} > \hat{\mu}_{p_3,a_1}$ and correctly returning m_{s1} .

While this richer action space of all matchings (regardless of stability) allows us to circumvent the shortcomings outlined in Example 1, it also permits and even encourages the selection of matchings that are known with high probability not to be stable. Although this approach efficiently identifies a true stable matching (and is indeed necessary, as the example shows), it might not be desirable in some applications to implement such "unfair" matchings as intermediate actions. In our case, we decide to focus on the pure exploration setting, aiming only to quickly identify the optimal stable matching.

3.3 Probably Correct Optimal Stable Matching

We introduce the notion of *Probably Correct Optimal Stable Matching* (PCOS). This concept is similar to the 'Probably Approximately Correct' (PAC) setting [7], where the goal is to find an ϵ -optimal arm with high probability using as few samples as possible in the context of MAB. Note that in the context of stable matchings it is non-trivial to define what an *approximation* of the optimal stable matching is. Consider Example 1 and assume for player p_3 the gap between true mean rewards $\Delta_{p_3,a_1,a_3} = |\mu_{p_3,a_1} - \mu_{p_3,a_3}|$ is less than ϵ . Then even an ϵ -approximation of the true means is not sufficient to identify the true optimal stable matching m_{s1} and instead m_{s2} is returned. In case one is tempted to consider approximations w.r.t. the sum of players' rewards over their matches we remark that here the approximation value of m_{s2} towards m_{s1} depends on the reward gaps Δ_{p_1,a_1,a_2} and Δ_{p_2,a_1,a_2} which might be arbitrarily large. We thus leave more extensive discussions around approximate solutions to future work and instead focus on the following solution concept.

DEFINITION 1 (PROBABLY CORRECT OPTIMAL STABLE MATCHING). We say that an algorithm A is a δ -PCOS algorithm with sample

complexity T , if it outputs the optimal stable matching, m_s^* , with probability at least $1 - \delta$ after at most T sampled matchings.

REMARK 1. We measure the sample complexity in terms of the number of matchings performed by algorithm A.

As correctly identifying the optimal stable matching is closely related to correctly identifying the agent's preferences we introduce the following two notions.

DEFINITION 2 (COMPLETELY CORRECT PREFERENCES). We say that the estimated preferences $\hat{\pi}$ of players are completely correct if they are exactly the same as the true preferences π , i.e., $\hat{\pi}[i] = \pi[i] \forall i \in [1, \dots, K]$.

DEFINITION 3 (PARTIALLY CORRECT PREFERENCES UP TO AN ARM). We say that the estimated preferences $\hat{\pi}$ of players are partially correct up to an arm $a \in \mathcal{A}$ if they are correct up to the position $r_p[a]$ with the true preference π , i.e., $\hat{\pi}[i] = \pi[i] \forall i \leq r_p[a]$.

We formally state the relation between the probability of the DA algorithm returning the optimal stable matching and the preference estimates being correct.

PROPOSITION 2. Let $m_{\hat{\pi}}$ be the output of the DA algorithm running using estimated preferences $\hat{\pi}$ for the agents \mathcal{A} . The probability that $m_{\hat{\pi}}$ is equal to the optimal stable matching m_s^* for the true preferences π is at least as high as the probability of those estimates $\hat{\pi}$ being correct, i.e.:

$$P(m_{\hat{\pi}} = m_s^*) \geq P(\hat{\pi}_i = \pi_i \forall i \in \mathcal{A}).$$

PROOF. This follows directly from the fact that the Deferred Acceptance algorithm with correct preferences always returns the optimal stable matching m_s^* . \square

4 NAIVE UNIFORM EXPLORATION

We begin our analysis by considering an algorithm that uniformly explores every available pair of agents, similarly to the ETC algorithm in the regret minimization setting [16]. Although ETC can achieve sublinear player-optimal stable regret, this does not always imply convergence to the correct optimal stable matching. Here, we provide a uniform sampling strategy that can identify the true optimal stable matching with high probability.

Our Naive Uniform Exploration (NUE) Algorithm uniformly samples matchings such that each pair of agents is sampled for a fixed number times, depending on the minimum reward difference between the arms. After the exploration rounds, the algorithm estimates the preferences of the players using the sample mean of the rewards for the arms and commits to the matching produced by the DA algorithm using these estimated preferences.

THEOREM 3. Let $\Delta_{\min} = \min_{p \in \mathcal{P}, i, j \in \mathcal{A}} \Delta_{p,i,j}$. Algorithm 1 is a δ -PCOS algorithm, and the number of matchings is bounded by:

$$O(K \frac{\ln(KN/\delta)}{\Delta_{\min}^2}). \quad (1)$$

PROOF SKETCH. The sample complexity follows from the definition of the algorithm, as $O(Kh) = O(K \ln(KN/\delta)/\Delta_{\min}^2)$.

We now prove that Algorithm 1 is an δ -PCOS algorithm, i.e., $\Pr(m_{\hat{\pi}} = m_s^*) \leq 1 - \delta$. Proposition 2 implies that $\Pr(m_{\hat{\pi}} \neq m_s^*) \leq$

Algorithm 1 Naive Uniform Exploration (NUE)

Require: $\delta > 0, \mathcal{P}, \mathcal{A}, \Delta_{\min} = \min_{p \in \mathcal{P}, i, j \in \mathcal{A}} \Delta_{p,i,j}, \{\pi_a\}_{a \in \mathcal{A}}$

- 1: $h = \lceil \frac{2 \ln(2KN/\delta)}{\Delta_{\min}^2} \rceil$
- 2: **for** t in $\{1, \dots, hK\}$ **do**
- 3: $m(p_i) = a_j, j = (t + i - 2 \bmod K) + 1 \forall i \in \{1, \dots, N\}$
- 4: Sample m and update $\hat{\mu}_{p_i, m(p_i)}(t) \forall i \in \{1, \dots, N\}$
- 5: **end for**
- 6: $\hat{\pi}_p = \arg \text{sort}_{a \in \mathcal{A}} \hat{\mu}_{p,a} \forall p \in \mathcal{P}$
- 7: $m_{\hat{\pi}} = DA(\{\hat{\pi}_p\}_{p \in \mathcal{P}}, \{\pi_a\}_{a \in \mathcal{A}})$
- 8: **return** $m_{\hat{\pi}}$

$\Pr(\bigcup_{p \in \mathcal{P}} \hat{\pi}_p \neq \pi_p)$. Using the union bound over the set of players \mathcal{P} and arms \mathcal{A} we have $\Pr(m_{\hat{\pi}} \neq m_s^*) \leq \sum_{p \in \mathcal{P}} \sum_{i=1}^N \Pr(\hat{\pi}_p[i] \neq \pi_p[i])$. Thus it is sufficient to show $\Pr(\hat{\pi}_p[i] \neq \pi_p[i]) \leq \delta/NK$.

The event that $\hat{\pi}_p[i] \neq \pi_p[i]$ can only occur if the player p wrongly orders two consecutive arms a and a' , i.e., $\hat{\mu}_{p,a} \leq \hat{\mu}_{p,a'}$ when $\mu_{p,a} \geq \mu_{p,a'}$. By construction of the sampled matchings (see line 3 in Algorithm 1), every player p receives h many rewards for every arm. Similar to the proof of Theorem 6 in [7] and by using Hoeffding's inequality, we have that $\Pr(\hat{\mu}_{p,a} \leq \hat{\mu}_{p,a'}) \leq \Pr(\hat{\mu}_{p,a} \leq \mu_{p,a} - \Delta_{p,a,a'}/2) + \Pr(\hat{\mu}_{p,a'} \geq \mu_{p,a'} + \Delta_{p,a,a'}/2) \leq 2 \exp(-\frac{\Delta_{\min}^2}{2} h) \leq \delta/NK$ which concludes the proof. \square

5 ELIMINATION ALGORITHM

In this section, we propose an elimination algorithm, similar to the one in [7] for MAB, which improves sample complexity and does not require prior knowledge of the reward differences $\Delta_{p,a,a'}$ of the agents. The key idea of our Elimination Algorithm (Algorithm 2) is to successively eliminate player-arm pairs (p, a) when the position of arm a in player p 's preference list can be determined with high probability. Eliminating every player-arm pair can eventually guarantee the correct estimation of the preferences of players with high probability, and thus identification of the optimal stable matching, as implied by Proposition 2.

Algorithm 2 operates in rounds, where in each round t , we sample matchings such that each player $p \in \mathcal{P}$ receives a reward from every non-eliminated arm from the set of available arms S_p , exactly ones (line 5). For this we compute a *minimal matching cover*, i.e., a cardinality-minimal set of matchings \mathcal{M} that cover all remaining pairs $\{(p, a) \in \mathcal{P} \times \mathcal{A} \mid a \in S_p\} \subseteq \bigcup_{m \in \mathcal{M}} m$. We denote $\mathcal{M}(X)$ to be a minimal matching cover on a given set of pairs X .

In the case of a bipartite graph, finding a minimal matching cover can be framed as a minimum edge coloring problem, where the edges of one color correspond to a matching in the matching cover [8]. According to the König-Hall Theorem, the optimal number of colors required is equal to the maximum degree of the graph – in our case the maximal degree $\deg_t = \Delta(G(E_t))$ of the bipartite graph $G(E_t)$ with edges E_t corresponding to available pairs $\{(p, a) \in \mathcal{P} \times \mathcal{A} \mid a \in S_p\}$. Thus, in every round t , we sample \deg_t matchings (see also Appendix C in the full version of our paper [1]).

The algorithm terminates once we eliminate all arms for every player, i.e., $S_p = \emptyset$ for all $p \in \mathcal{P}$. We use the elimination rule that eliminates a (p, a) pair when the arm a has no overlapping confidence interval $C_{p,a}$ with other available arms in S_p for player p .

Together with the use of anytime confidence intervals¹, the elimination rule can guarantee the position of arm a in the preference list of player p with high probability.

Algorithm 2 Elimination Algorithm

Require: $\delta > 0, \mathcal{P}, \mathcal{A}, \{\pi_a\}_{a \in \mathcal{A}}$

```

1:  $S_p = \mathcal{A} \ \forall p \in \mathcal{P}$ 
2:  $t = 1$ 
3:  $\hat{\mu}_{p,a}(t) = 0 \ \forall p \in \mathcal{P}, a \in \mathcal{A}$ 
4: while  $|\cup_{p \in \mathcal{P}} S_p| \geq 1$  do
5:   Sample all matchings  $m \in \mathfrak{M}(\{(p, a) \in \mathcal{P} \times \mathcal{A} \mid a \in S_p\})$ 
6:   Update  $\hat{\mu}_{p,a}(t) \ \forall p \in \mathcal{P}$  and  $a \in S_p$ 
7:    $B_t = \sqrt{\frac{\ln(4KNt^2/\delta)}{2t}}$ 
8:    $C_{p,a} = [\hat{\mu}_{p,a}(t) \pm B_t] \ \forall p \in \mathcal{P}$  and  $a \in S_p$ 
9:    $S_p = S_p \setminus \{a : C_{p,a} \cap C_{p,j} = \emptyset \ \forall j \in \mathcal{A} \setminus \{a\}\} \ \forall p \in \mathcal{P}$ 
10:   $t = t + 1$ 
11: end while
12:  $\hat{\pi}_p = \arg \text{sort}_{a \in \mathcal{A}} \hat{\mu}_{p,a} \ \forall p \in \mathcal{P}$ 
13: return  $DA(\{\hat{\pi}_p\}_{p \in \mathcal{P}}, \{\pi_a\}_{a \in \mathcal{A}})$ 

```

As a starting point for the overall sample complexity of Algorithm 2, we first calculate how many samples $t_{p,a}$ that are sufficient to eliminate an arm a from a player p 's list of available arms S_p as in line 9 of the algorithm.

LEMMA 1. For a player $p \in \mathcal{P}$ and arm $a \in \mathcal{A}$, let $\Delta_{p,a} = \min_{a' \in \mathcal{A} \setminus \{a\}} \Delta_{p,a,a'}$. With probability at least $1 - \delta$, the number of samples $t_{p,a}$ needed to eliminate an arm a from S_p in line 9 of Algorithm 2, is at most

$$t_{p,a} = O\left(\frac{\ln(KN/\delta\Delta_{p,a})}{\Delta_{p,a}^2}\right). \quad (2)$$

PROOF SKETCH. Let \mathcal{E} denote the event that, for all time steps t the expected rewards $\mu_{p,a}$ lie in the confidence interval $CI_{p,a}(t)$ for every pair of player p and arm a i.e. $\mu_{p,a} - B_t \leq \hat{\mu}_{p,a}(t) \leq \mu_{p,a} + B_t \ \forall t, a \in \mathcal{A}, p \in \mathcal{P}$. Using Hoeffding's inequality, we can show that the values $B_t = \sqrt{\frac{\ln(4KNt^2/\delta)}{2t}}$ in Algorithm 2 correspond to bounds of any-time confidence intervals $CI_{p,a}(t) = [\hat{\mu}_{p,a}(t) - B_t, \hat{\mu}_{p,a}(t) + B_t]$. Consequently, \mathcal{E} is true w.p.a. $1 - \delta$.

Now assume \mathcal{E} is true. For a player p consider two arms a, a' with $\mu_{p,a} > \mu_{p,a'}$. Under the event \mathcal{E} and for t such that $\Delta_{p,a,a'} > 4B_t$, the arms will have no overlapping confidence intervals as $(\hat{\mu}_{p,a} - B_t) - (\hat{\mu}_{p,a'} + B_t) \geq \mu_{p,a} - \mu_{p,a'} - 4B_t > 0$. So in order for an arm a to have no overlapping confidence intervals with any other arm we need $\Delta_{p,a} = \min_{a' \in \mathcal{A} \setminus \{a\}} \Delta_{p,a,a'} > 4B_t$, which holds for $t_{p,a} = O\left(\ln(KN/\delta\Delta_{p,a})/\Delta_{p,a}^2\right)$ using a similar analysis with [7].

Consequently, after $t_{p,a}$ samples of all arms, no confidence interval $CI_{p,a'}(t) \ \forall a' \in \mathcal{A} \setminus \{a\}$ overlaps with $CI_{p,a}(t)$ and a is eliminated from S_p , with probability at least $1 - \delta$. \square

We can now state the sample complexity and show Algorithm 2 is a δ -PCOS algorithm, with a detailed proof in Appendix D of the full version of our paper [1].

¹Anytime confidence intervals ensure valid coverage probabilities at any time t .

THEOREM 4. Let $\mathcal{P}(S) = \{r \subseteq S : \Delta(G(S)) - \Delta(G(S \setminus r)) = 1\}$ denote the set of pairs that we have to eliminate to reduce the degree of the graph $G(S)$ with $S = \{(p, a) \in \mathcal{P} \times \mathcal{A}\}$. Let also $r_0 = \emptyset$ and

$$r_i = \arg \min_{r \in \mathcal{P}(S \setminus \bigcup_{j < i} r_j)} \max_{(p,a) \in r} t_{p,a} \ \forall i = 1, \dots, K$$

Algorithm 2 is a δ -PCOS algorithm, and with probability at least $1 - \delta$, the number of matching samples is bounded by:

$$O\left(\sum_{s=1}^K \max_{(p,a) \in r_s} t_{p,a}\right) \quad (3)$$

PROOF SKETCH. Using Hoeffding's inequality, we can show that the values B_t (in Algorithm2) define any-time confidence intervals, i.e., the event \mathcal{E} , where at any time-step t the expected rewards $\mu_{p,a}$ lie in the confidence interval $C_{p,a}(t)$ for every pair of player p and arm a , holds with probability at least $1 - \delta$.

Under the event \mathcal{E} , once two confidence intervals are not overlapping for some arms a, a' , we can determine their relative order in player p 's preferences. This still holds true even when taking more samples of one of the arms, as the confidence intervals only shrinks. We can thus eliminate the pair (p, a) once the order of the arm a towards all other arms can be determined with high probability. Thus, upon termination of the algorithm, we will have determined the correct preferences $\hat{\pi}_p = \pi_p \ \forall p \in \mathcal{P}$ and by Proposition 2, Algorithm 2 will output m_s^* with w.p.a. $1 - \delta$.

For the sample complexity, consider that at any time t we sample every available pair, $\{(p, a) \in \mathcal{P} \times \mathcal{A} \mid a \in S_p\}$, using the matchings from a minimal matching cover of the bipartite graph with edges corresponding to currently available pairs (see line 5). Such a minimal matching cover has \deg_t many matchings. Thus, once the maximal degree is reduced, less matchings have to be sampled.

We thus consider $s = 1, \dots, K$ phases where each phase corresponds to the elimination of the subset of pairs r_s to reduce the degree of the graph. So under the event \mathcal{E} , we can use Lemma 1 to define r_s as the set of pairs with the lowest sample complexity i.e:

$$r_i = \arg \min_{r \in \mathcal{P}(S \setminus \bigcup_{j < i} r_j)} \max_{(p,a) \in r} t_{p,a} \ \forall i = 1, \dots, K$$

with

$$\mathcal{P}(S) = \{r \subseteq S : \Delta(G(S)) - \Delta(G(S \setminus r)) = 1\}$$

denote the subsets of pairs that can reduce the degree of the graph. Note that here the r_1, \dots, r_K form a partition of all player-arm pairs.

In every phase s , we perform $t_s - t_{s-1}$ iterations, where t_s are sufficient number of steps to eliminate r_s determined by the player-arm pair $(p, a) \in r_s$ that takes the longest to eliminate i.e: $t_s = \max_{(p,a) \in r_s} t_{p,a}$ with $t_0 = 0$.

Finally, in each step t , every matching cover consists of $K - s + 1$ matchings. Thus, under the event \mathcal{E} which holds w.p.a. $1 - \delta$, the total number of matching samples is given by:

$$\sum_{s=1}^K (K - s + 1)(t_s - t_{s-1}) = \sum_{s=1}^K t_s = O\left(\sum_{s=1}^K \max_{(p,a) \in r_s} t_{p,a}\right)$$

\square

REMARK 2. Note, that in the case where $N = K$, the set of pairs r_s that we have to eliminate, corresponds to a perfect matching m_s as in this case the graph of available pairs is regular at each phase s .

REMARK 3. We can construct a worst-case instance, where Algorithm 2 uniformly samples all player-arm pairs. In particular, if the differences in the expected rewards are equal i.e. $\Delta_{p,a} = \Delta \forall p \in \mathcal{P}, a \in \mathcal{A}$, leading to sample complexity:

$$O\left(\sum_{s=1}^K \max_{(p,a) \in r_s} t_{p,a}\right) = O\left(K \frac{\ln(KN/\delta\Delta)}{\Delta^2}\right).$$

6 IMPROVED ELIMINATION ALGORITHM

In this section, we propose an improved version of the elimination algorithm, Algorithm 3, based on the observation that to identify the optimal stable matching, we only need to correctly estimate the preferences of the players up to the position of the matching partner in the optimal stable matching. We formalize this in the Lemma below.

LEMMA 2. Let m_s^* be the true optimal stable matching according to preferences π . The output, $m_{\hat{\pi}}$, of the DA algorithm using preferences $\hat{\pi}$ that are partially correct up to $m_{\hat{\pi}}(p)$ for every player p , i.e., $\hat{\pi}_p[i] = \pi_p[i] \forall i \leq r_p[m_{\hat{\pi}}(p)] \forall p \in \mathcal{P}$, is equal to the true optimal stable matching m_s^* .

PROOF. The DA algorithm sequentially executes proposals from players starting with their most preferred arm and never backtracks to a previously made proposal. Further, it halts once the player optimal stable matching has been found. Consequently it only considers the matching partner $m_{\hat{\pi}}(p)$ and all higher ranked arms in the optimal stable matching according to $\hat{\pi}$. So if the preferences π_p are correct up to the positions of the $m_{\hat{\pi}}$ then is equal to the true optimal stable matching m_s^* . \square

We can modify the stopping criteria of the Elimination Algorithm of the previous section according to Lemma 2. Algorithm 3 terminates when it eliminates the arms up to the stable matching partner for every player. Specifically, after each round t , we calculate an estimate of the player optimal stable matching \hat{m}_t from the DA algorithm using the estimated preferences $\hat{\pi}$ from the sample means. The algorithm terminates if for all players $p \in \mathcal{P}$, every the stable matching partner $\hat{m}_t(p)$ and all higher ranked arms have been eliminated. As the algorithm proceeds the stable matching \hat{m}_t changes until we eventually can reach a state where our termination criteria holds.

THEOREM 5. Let $\mathcal{P}(S) = \{r \subseteq S : \Delta(G(S)) - \Delta(G(S \setminus r)) = 1\}$ denote the set of pairs that we have to eliminate to reduce the degree of the graph $G(S)$ with $S = \{(p, a) \in \mathcal{P} \times \mathcal{A}\}$, $r_0 = \emptyset$ and

$$r_i = \arg \min_{r \in \mathcal{P}(S \setminus \bigcup_{j < i} r_j)} \max_{(p,a) \in r} t_{p,a} \forall i = 1, \dots, K$$

Let also $t_{\max} = \max_{(p,a) \in \mathcal{P} \times \mathcal{A} : a \succeq_p m_s^*(p)} t_{p,a}$, and $n \in \{1, \dots, K\}$ be the index s.t. $\max_{(p,a) \in r_{n-1}} t_{p,a} \leq t_{\max} \leq \max_{(p,a) \in r_n} t_{p,a}$.

Algorithm 3 is a δ -PCOS algorithm, and with probability at least $1 - \delta$, the number of matching samples is bounded by:

$$O\left(\sum_{s=1}^{n-1} \max_{(p,a) \in r_s} t_{p,a} + (K - n + 1)t_{\max}\right) \quad (4)$$

PROOF. Algorithm 3 is essentially the same as Algorithm 2 with a different stopping criterion (see lines 3 and 11), and many of arguments from the proof of Theorem 4 transfer. In particular, at

Algorithm 3 Improved Elimination Algorithm

Require: $\delta > 0, \mathcal{P}, \mathcal{A}, \{\pi_a\}_{a \in \mathcal{A}}$

```

1:  $t = 1, S_p = \mathcal{A} \forall p \in \mathcal{P}, S = S_p$ 
2:  $\hat{\mu}_{p,a}(t) = 0 \forall p \in \mathcal{P}, a \in \mathcal{A}$ 
3: while  $|S| \geq 1$  do
4:   Sample all matchings  $m \in \mathfrak{M}(\{(p, a) \in \mathcal{P} \times \mathcal{A} \mid a \in S_p\})$ 
5:   Update  $\hat{\mu}_{p,a}(t) \forall p \in \mathcal{P}$  and  $a \in S_p$ 
6:    $B_t = \sqrt{\frac{\ln 4KNt^2/\delta}{2t}}$ 
7:    $C_{p,a} = [\hat{\mu}_{p,a}(t) \pm B_t] \forall a \in S_p, p \in \mathcal{P}$ 
8:    $S_p = S_p \setminus \{a : C_{p,a} \cap C_{p,j} = \emptyset \forall j \in \mathcal{A} \setminus \{a\}\}$ 
9:    $\hat{\pi}_p = \arg \text{sort}_{a \in \mathcal{A}} \hat{\mu}_{p,a} \forall p \in \mathcal{P}$ 
10:   $m_t = DA(\{\hat{\pi}_p\}_{p \in \mathcal{P}}, \{\pi_a\}_{a \in \mathcal{A}})$ 
11:   $S = \bigcup_{p \in \mathcal{P}} \{a \in S_p : a \succeq_{\hat{\pi}_p} m_t(p)\}$ 
12:   $t = t + 1$ 
13: end while
14: return  $m_t$ 
```

the time of elimination of an arm a from S_p , the arm can be correctly ordered w.r.t. all other arms in p 's preferences with high probability. Thus at the time of termination, with probability at least $1 - \delta$, the condition of Lemma 2 is satisfied and Algorithm 3 outputs the correct optimal stable matching m_s^* .

For the sample complexity, we can again consider phases s where we eliminate a subset of pairs r_s after at most $t_s = \max_{(p,a) \in r_s} t_{p,a}$ number of samples, under the event \mathcal{E} . However, with the stopping criterion in Algorithm 3, we might not have to eliminate every r_1, \dots, r_K before stopping. In fact, by Lemma 2 we only need to eliminate the optimal matching partner and all higher ranked arms for every player i.e. the pairs $\{(p, a) \in \mathcal{P} \times \mathcal{A} : a \succeq_p m_s^*(p)\}$. For this, a maximal number of samples t_{\max} from every pair are needed.

Now, consider the last phase n in which all arms $a \succeq_p m_s^*(p)$ are either eliminated from S_p or sampled t_{\max} many times, i.e., $t_{n-1} \leq t_{\max} \leq t_n$. The algorithm terminates within phase n , with a total number of matching: $\sum_{s=1}^{n-1} (K - s + 1)(t_s - t_{s-1}) + (K - n + 1)(t_{\max} - t_{n-1}) = \sum_{s=1}^{n-1} (K - s + 1)t_s + (K - n + 1)t_{\max}$. To conclude, under the event \mathcal{E} which holds w.p.a. $1 - \delta$, the algorithm terminates with the total number of matchings bounded by:

$$O\left(\sum_{s=1}^{n-1} \max_{(p,a) \in r_s} t_{p,a} + (K - n + 1)t_{\max}\right) \quad \square$$

REMARK 4. We can construct a worst-case instance, where Algorithm 3 uniformly samples all player-arm pairs. In particular, if the differences in the expected rewards are equal i.e. $\Delta_{p,a} = \Delta \forall p \in \mathcal{P}, a \in \mathcal{A}$, leading to sample complexity:

$$O\left(K \frac{\ln(KN/\delta\Delta)}{\Delta^2}\right).$$

As for Algorithm 2, Algorithm 3 will, in the worst case, have the same sample complexity as a strategy that uniformly samples pairs until they are eliminated, e.g., when all expected reward gaps are equal, i.e., $\Delta_{p,a} = \Delta \forall p \in \mathcal{P}, a \in \mathcal{A}$ (see also Remark 3). However, in practice, we expect Algorithm 3 to perform much better and terminate earlier, particularly when the reward gaps are larger for higher-ranked consecutive arm pairs.

7 ADAPTIVE SAMPLING ALGORITHM

The algorithms introduced in the previous sections uniformly sample the arms for each player, until their confidence intervals are sufficiently separated. Here we propose an approach to adaptively sample player-arm pairs at every round.

Algorithm 4, leverages insights from Lemma 2 to dynamically define the set of arms requiring further exploration, denoted as S_p . Specifically, in each round t , we estimate the optimal stable matching \hat{m}_t from the empirical preferences $\hat{\pi}$. According to Lemma 2, for this matching to be correct, the arms must be accurately ranked for the players up to the partner in the optimal stable matching. Consequently, in each round t and for each player p , we sample the set of arms in order to distinguish the confidence intervals of the arms up to the stable matching partner $m_t(p)$ (see line 11). Due to the adaptive selection of arms, each arm has a distinct confidence margin $B_{p,a}$ (line 5), which varies based on the number of times we sample a player-arm pair $t_{p,a}$.

Algorithm 4 Adaptive Sampling Algorithm

Require: $\delta > 0, \mathcal{P}, \mathcal{A}, \{\pi_a\}_{a \in \mathcal{A}}$

```

1:  $t_{p,a} = 0, S_p = \mathcal{A} \forall p \in \mathcal{P}$ 
2:  $\hat{\mu}_{p,a}(t) = 0 \forall p \in \mathcal{P}, a \in \mathcal{A}$ 
3: while  $|\cup_{p \in \mathcal{P}} S_p| \geq 1$  do
4:   Sample all matchings  $m \in \mathfrak{M}(\{(p,a) \in \mathcal{P} \times \mathcal{A} \mid a \in S_p\})$ 
5:   Update  $\hat{\mu}_{p,a}(t)$  and  $t_{p,a} \forall p \in \mathcal{P}$  and  $a \in S_p$ 
6:    $B_{p,a} = \sqrt{\frac{\ln 4KN t_{p,a}^2 / \delta}{2t_{p,a}}} \forall a \in \mathcal{A}, p \in \mathcal{P}$ 
7:    $C_{p,a} = [\hat{\mu}_{p,a} \pm B_{p,a}] \forall a \in \mathcal{A}, p \in \mathcal{P}$ 
8:    $\hat{\pi}_p = \arg \text{sort}_{a \in \mathcal{A}} \hat{\mu}_{p,a} \forall p \in \mathcal{P}$ 
9:    $m_t = DA(\{\hat{\pi}_p\}_{p \in \mathcal{P}}, \{\pi_a\}_{a \in \mathcal{A}})$ 
10:   $A[p] = \{a \in \mathcal{A} : a \succeq_{\hat{\pi}_p} m_t(p)\} \forall p \in \mathcal{P}$ 
11:   $S_p = \{a \in \mathcal{A} \mid \exists a' \in \mathcal{A} \setminus \{a\} : C_{p,a'} \cap C_{p,a} \neq \emptyset \text{ and } \{a, a'\} \cap A[p] \neq \emptyset\} \forall p \in \mathcal{P}$ 
12:   $t = t + 1$ 
13: end while
14: return  $m_t$ 

```

THEOREM 6. *Algorithm 4 is a δ -PCOS algorithm.*

PROOF. First, note that, as shown for Algorithm 2 and 3, the $B_{p,a}$ define anytime confidence intervals, i.e., the event \mathcal{E} — where the expected reward lies within the confidence intervals at any time t for all pairs — holds with probability at least $1 - \delta$.

So under \mathcal{E} , consider a round where the estimate of the stable matching m_t is incorrect, which occurs only when at least one player p has wrong preferences $\hat{\pi}_p$ up to the stable match. This implies that at least one arm a with $a \succeq_{\hat{\pi}_p} m_t(p)$ has an overlapping confidence interval with another arm a' , defining the active set of arms in round t , e.g., $S_p = \{a, a'\}$. In this situation, the algorithm samples the arms in the active set S_p until their confidence intervals are sufficiently separated or until we find ourselves with a different matching $m_{t'}$ and a new active set of arms. Since our confidence intervals only shrink, our algorithm will eventually terminate with the correct preferences up to the stable match, and thus with the correct optimal stable matching with probability at least $1 - \delta$. \square

The adaptive selection of agent pairs complicates the analysis of sample complexity, which we leave for future work. Instead, we test the algorithms in practice in the following section.

8 SIMULATIONS

In this section, we perform simulations on random instances to further evaluate the performance of our algorithms for fixed $\delta = 0.1$. The rewards for each pair are drawn from a Bernoulli distribution, i.e., $X_{p,a} \sim \text{Bern}(\mu_{p,a})$, with $\mu_{p,a} = \mathbb{E}[X_{p,a}]$. We perform different experiments by varying the number of agents with $N = K$, while we explore two different reward settings: (1) random expected rewards, and (2) random expected rewards with decreasing gaps. For each experiment, we generate 100 random instances. The code for this work is available at <https://github.com/a-athanasopoulos/PACOS>.

In both reward settings, we randomly generate the preferences of the agents π for each instance. To create the respective expected rewards μ for the pairs, we first sample $K - 1$ arbitrary reward gaps $\Delta_{p,i,i+1}$ for each player p and $i = 1, \dots, K - 1$ from a Dirichlet distribution $\text{Dir}(\alpha = 1)$, while we set $\Delta_{p,0,1} = 0$. For computational reasons, we normalize these values $\Delta_{p,i,i+1}$ to ensure that none exceed 0.05. In **Reward Setting 1**, we set the expected rewards to $\mu_{p,\pi_p[i]} = \sum_{j \leq K-i} \Delta_{p,j,j+1}$, while in **Reward Setting 2**, we first sort the reward gaps in increasing order. The second setting ensures reduced sample complexity for the algorithms that employs the stopping rule according to our Lemma 2, as the preferences up to the stable match are easier to distinguish, i.e., $\Delta_{p,\pi[i],\pi[i+1]} \geq \Delta_{p,\pi[i+1],\pi[i+2]} \forall i = 1, \dots, K - 1$.

We compare the (1) **Elimination Algorithm**, (2) **Improved Elimination Algorithm**, and (3) **Adaptive Sampling Strategy**. In addition, we consider a variant of the NUE algorithm, the (4) **Uniform Sampling Strategy** that uniformly samples every pair

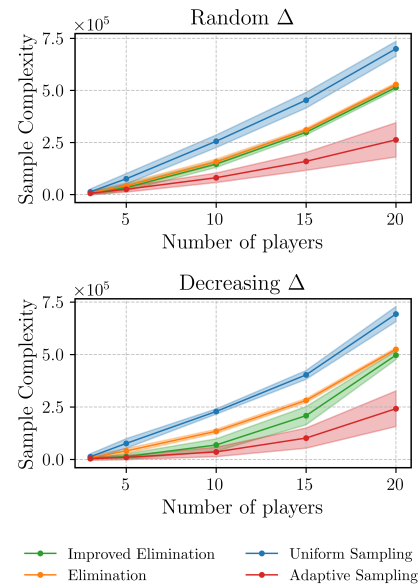


Figure 1: Sample complexity for the proposed algorithms for the two different reward settings, averaged over the runs.

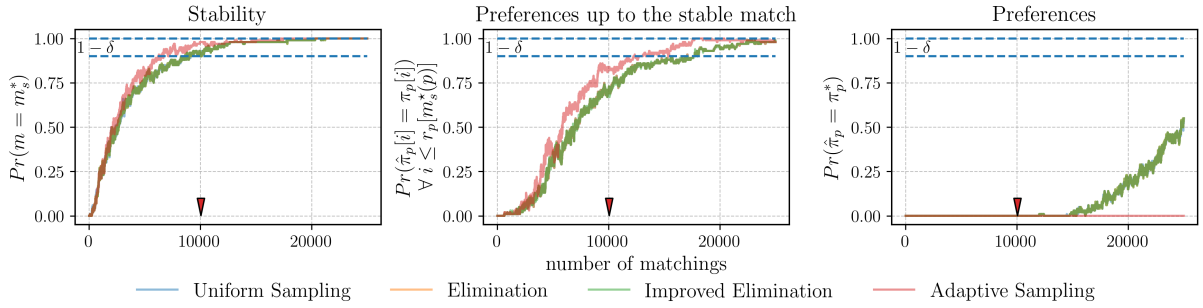


Figure 2: Any-time performance of the algorithm for the first instance with 20 agents on each side. The figure illustrates the average number of times the algorithms are able to identify (left) the optimal stable matching, (middle) the correct preferences up to the stable match for every player, and (right) the correct preferences for every player, after each matching.

until there are no overlapping confidence intervals, as the algorithm introduced in Section 4 requires knowledge of the minimum reward difference Δ and has a fixed sample complexity.

8.1 Sample Complexity

First note, that in every experiment the algorithms always return the correct stable matching, similar to the study on MAB [13]. We further discuss the anytime performance of the algorithms in Section 8.2. In Figure 1, we present the average and the standard deviation of the sample complexity over the instances, for both preference settings, respectively.

In the first setting, we observe that the elimination algorithms behave similarly. This is because the randomly generated preferences do not allow the Improved Elimination Algorithm to terminate early. In the second setting, however, the Improved Elimination Algorithm outperforms the standard Elimination strategy. Additionally, the Uniform Sampling Strategy requires more samples, even with 20 players, where Remark 3 indicates similar sample complexity, as the differences in expected rewards are similar ($\Delta \approx 0.05$). This occurs because our theoretical sample complexity measures the sufficient number of samples, while in practice some arms can be eliminated earlier. Finally, the Adaptive Sampling Strategy outperforms all other algorithms in both settings, as it dynamically refines exploration based on the agents' preferences.

8.2 Anytime Performance

Now we study the performance of the algorithms at each time step t , similar to the approaches described in [2, 13] for the MAB. More specifically, after each time step t , we can check whether the DA algorithm using the current preferences, outputs the correct m_s^* . We also checked if the preferences are completely correct, and correct up to the stable match, respectively. In Figure 2, show the results averaged over 100 runs for Reward Setting 1 with 20 players. The results for the remaining experiments can be found in Appendix E of the full version of our paper [1].

First, note that the curves for the two elimination algorithms overlap, as the only difference among them lies in the stopping criterion. The same is also true for the uniform sampling strategy. This is because elimination begins after the plotted time window (after 54,000 matchings), so the algorithms uniformly sample arms

until that point. On the other hand, the adaptive algorithm quickly reduces the number of matchings to be explored, as indicated by the red arrow, which marks the average time when we need to explore less agents. This allows the algorithm to efficiently identify the preferences up to the stable pair and, consequently, the optimal stable matching. In addition, the algorithm fails to accurately estimate the preferences, as it focuses only on exploring the arms crucial for stability. Finally, we can also observe that the probability of achieving the optimal stable matching is greater than the other metrics, as indicated by Proposition 2 and Lemma 2.

9 CONCLUSION & FUTURE WORK

In this work, we consider the stable marriage problem under uncertain preferences on one side of the market. Our objective was to develop algorithms that efficiently identify the true optimal stable matching with high probability. To this end, we proposed the novel concept of a Probably Correct Optimal Stable Matching and present several algorithms with theoretical analyses of correctness and sample complexity. Finally, we support our theoretical results with an empirical evaluation of our algorithms' performance in practice.

There are several interesting directions for future research. One promising area is the hardness analysis, particularly establishing a lower bound on the sample complexity for the pure exploration problem. Additionally, analyzing the principles of the Elimination Algorithm in the regret minimization setting [16] and examining its relationship to the pure exploration problem is an important direction. Another straightforward extension is to investigate how these algorithms perform when both sides of the market are uncertain about their preferences (discussed in Appendix A [1]). The sample complexity of the Adaptive Sampling Strategy remains an open question, while experiments showcase its superiority over the other strategies. Finally, exploring alternative solution concepts, such as popular matchings [11] and median stable matchings [23], or studying alternative models for two-sided markets [21], can potentially broaden the applicability of our approach.

ACKNOWLEDGMENTS

This project was partially supported by the Norwegian Research Council under project 302203 "Algorithms and Models for Socially Beneficial AI".

REFERENCES

- [1] Andreas Athanasiopoulos, Anne-Marie George, and Christos Dimitrakakis. 2025. Probably Correct Optimal Stable Matching for Two-Sided Markets Under Uncertainty. *arXiv:2501.03018* [cs.LG]
- [2] P. Auer. 2002. Finite-time Analysis of the Multiarmed Bandit Problem.
- [3] Soumya Basu, Karthik Abinav Sankararaman, and Abishek Sankararaman. 2021. Beyond $\log^2(T)$ regret for decentralized bandits in matching markets. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, Virtual, 705–715. <https://proceedings.mlr.press/v139/basu21a.html>
- [4] Sarah H. Cen and Devavrat Shah. 2022. Regret, stability & fairness in matching markets with bandit learners. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics (Proceedings of Machine Learning Research, Vol. 151)*, Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (Eds.). PMLR, Virtual, 8938–8968. <https://proceedings.mlr.press/v151/cen22a.html>
- [5] Sanmay Das and Emir Kamenica. 2005. Two-sided bandits and the dating market. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence (Edinburgh, Scotland) (IJCAI'05)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 947–952.
- [6] Eyal Even-Dar, Shie Mannor, and Yishay Mansour. 2002. PAC Bounds for Multi-armed Bandit and Markov Decision Processes. In *Proceedings of the 15th Annual Conference on Computational Learning Theory (COLT '02)*. Springer-Verlag, Berlin, Heidelberg, 255–270.
- [7] Eyal Even-Dar, Shie Mannor, and Yishay Mansour. 2006. Action Elimination and Stopping Conditions for the Multi-Armed Bandit and Reinforcement Learning Problems. *J. Mach. Learn. Res.* 7 (dec 2006), 1079–1105.
- [8] Harold N. Gabow and Oded Kariv. 1978. Algorithms for edge coloring bipartite graphs. In *Proceedings of the Tenth Annual ACM Symposium on Theory of Computing (San Diego, California, USA) (STOC '78)*. Association for Computing Machinery, New York, NY, USA, 184–192. <https://doi.org/10.1145/800133.804346>
- [9] D. Gale and L. S. Shapley. 1962. College Admissions and the Stability of Marriage. *The American Mathematical Monthly* 69, 1 (1962), 9–15. <http://www.jstor.org/stable/2312726>
- [10] Guillaume Haeringer and Myrna Wooders. 2011. Decentralized Job Matching. *International Journal of Game Theory* 40 (02 2011), 1–28. <https://doi.org/10.1007/s00182-009-0218-x>
- [11] Chien-Chung Huang and Telikepalli Kavitha. 2013. Popular matchings in the stable marriage problem. *Information and Computation* 222 (2013), 180–194.
- [12] Meena Jagadeesan, Alexander Wei, Yixin Wang, Michael I Jordan, and Jacob Steinhardt. 2023. Learning equilibria in matching markets with bandit feedback. *J. ACM* 70, 3 (2023), 1–46.
- [13] Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. 2014. *lil' UCB*: An Optimal Exploration Algorithm for Multi-Armed Bandits. In *Proceedings of The 27th Conference on Learning Theory (Proceedings of Machine Learning Research, Vol. 35)*, Maria Florina Balcan, Vitaly Feldman, and Csaba Szepesvári (Eds.). PMLR, Barcelona, Spain, 423–439. <https://proceedings.mlr.press/v35/jamieson14.html>
- [14] Fang Kong, Junming Yin, and Shuai Li. 2022. Thompson Sampling for Bandit Learning in Matching Markets. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, Lud De Raedt (Ed.). International Joint Conferences on Artificial Intelligence Organization, Vienna, Austria, 3164–3170. <https://doi.org/10.24963/ijcai.2022/439> Main Track.
- [15] Tor Lattimore and Csaba Szepesvári. 2020. *Bandit algorithms*. Cambridge University Press, Cambridge, UK.
- [16] Lydia T. Liu, Horia Mania, and Michael I. Jordan. 2020. Competing Bandits in Matching Markets. *arXiv:1906.05363* [cs.LG]
- [17] Lydia T. Liu, Feng Ruan, Horia Mania, and Michael I. Jordan. 2021. Bandit Learning in Decentralized Matching Markets. *arXiv:2012.07348* [cs.LG]
- [18] David F. Manlove. 2013. *Algorithmics of Matching Under Preferences*. WORLD SCIENTIFIC, Singapore. <https://doi.org/10.1142/8591> *arXiv:https://www.worldscientific.com/doi/pdf/10.1142/8591*
- [19] Shie Mannor and John N. Tsitsiklis. 2004. The sample complexity of exploration in the multi-armed bandit problem. *Journal of Machine Learning Research* 5, Jun (2004), 623–648.
- [20] Alvin E. Roth. 1984. The evolution of the labor market for medical interns and residents: a case study in game theory. *Journal of political Economy* 92, 6 (1984), 991–1016.
- [21] Alvin E. Roth and Marilda A. Oliveira Sotomayor. 1990. *Two-Sided Matching: A Study in Game-Theoretic Modeling and Analysis*. Cambridge University Press, Cambridge, UK.
- [22] Abishek Sankararaman, Soumya Basu, and Karthik Abinav Sankararaman. 2021. Dominate or Delete: Decentralized Competing Bandits in Serial Dictatorship. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics (Proceedings of Machine Learning Research, Vol. 130)*, Arindam Banerjee and Kenji Fukumizu (Eds.). PMLR, Virtual, 1252–1260. <https://proceedings.mlr.press/v130/sankararaman21a.html>
- [23] Chung-Piaw Teo and Jay Sethuraman. 1998. The geometry of fractional stable matchings and its applications. *Mathematics of Operations Research* 23, 4 (1998), 874–891.
- [24] William R. Thompson. 1933. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 25, 3-4 (1933), 285–294.