



# Learning Equilibria in Matching Markets with Bandit Feedback

MEENA JAGADEESAN\*, ALEXANDER WEI\*, and YIXIN WANG, UC Berkeley EECS, USA  
MICHAEL I. JORDAN, UC Berkeley EECS and Statistics, USA  
JACOB STEINHARDT, UC Berkeley Statistics, USA

19

Large-scale, two-sided matching platforms must find market outcomes that align with user preferences while simultaneously learning these preferences from data. Classical notions of stability (Gale and Shapley, 1962; Shapley and Shubik, 1971) are, unfortunately, of limited value in the learning setting, given that preferences are inherently uncertain and destabilizing while they are being learned. To bridge this gap, we develop a framework and algorithms for learning stable market outcomes under uncertainty. Our primary setting is matching with transferable utilities, where the platform both matches agents and sets monetary transfers between them. We design an incentive-aware learning objective that captures the distance of a market outcome from equilibrium. Using this objective, we analyze the complexity of learning as a function of preference structure, casting learning as a stochastic multi-armed bandit problem. Algorithmically, we show that “optimism in the face of uncertainty,” the principle underlying many bandit algorithms, applies to a primal-dual formulation of matching with transfers and leads to near-optimal regret bounds. Our work takes a first step toward elucidating when and how stable matchings arise in large, data-driven marketplaces.

CCS Concepts: • **Theory of computation** → **Regret bounds; Exact and approximate computation of equilibria;**

Additional Key Words and Phrases: Multi-armed bandits, stable matchings, learning equilibria, preference structure

## ACM Reference format:

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, Article 19 (May 2023), 46 pages.  
<https://doi.org/10.1145/3583681>

## 1 INTRODUCTION

Data-driven marketplaces face the simultaneous challenges of learning agent preferences and aligning market outcomes with the incentives induced by these preferences. Consider, for instance,

\*M. Jagadeesan and A. Wei Equal contribution.

Meena Jagadeesan acknowledges support from the Paul and Daisy Soros Fellowship; Alexander Wei acknowledges support from an NSF Graduate Research Fellowship under Grant No. 1752814; and Michael Jordan acknowledges the Vannevar Bush Faculty Fellowship program under grant number N00014-21-1-2941.

Authors' addresses: M. Jagadeesan, A. Wei, M. I. Jordan, and J. Steinhardt, UC Berkeley EECS, 2121 Berkeley Way #47, Berkeley, CA 94720-1660; emails: [mjagadeesan@berkeley.edu](mailto:mjagadeesan@berkeley.edu), [awei@berkeley.edu](mailto:awei@berkeley.edu), [jordan@cs.berkeley.edu](mailto:jordan@cs.berkeley.edu), [jsteinhardt@berkeley.edu](mailto:jsteinhardt@berkeley.edu); Y. Wang, UC Berkeley EECS, 1085 S University Avenue, Ann Arbor, MI, 48109; email: [ywang@eecs.berkeley.edu](mailto:ywang@eecs.berkeley.edu).

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2023 Copyright held by the owner/author(s).

0004-5411/2023/05-ART19 \$15.00

<https://doi.org/10.1145/3583681>

online platforms that match two sides of a market to each other (e.g., Lyft, TaskRabbit, and Airbnb). On these platforms, customers are matched to service providers and pay for the service they receive. If agents on either side are not offered desirable matches at fair prices, then they would have an incentive to leave the platform and switch to a competing platform. Agent preferences, however, are often unknown to the platform and must be learned. When faced with uncertainty about agent preferences (and thus incentives), *when can a marketplace efficiently explore and learn market outcomes that align with agent incentives?*

We center our investigation around a model called *matching with transferable utilities*, proposed by Shapley and Shubik [42]. In this model, there is a two-sided market of customers and service providers. Each customer has a utility that they derive from being matched to a given provider and vice versa. The platform selects a matching between the two sides and assigns a monetary transfer between each pair of matched agents. Transfers are a salient feature of most real-world matching markets: riders pay drivers on Lyft, clients pay freelancers on TaskRabbit, and guests pay hosts on Airbnb. An agent's net utility is their value for being matched to their partner plus the value of their transfer (either of which can be negative in the cases of costs and payments). In matching markets, the notion of *stability* captures alignment of a market outcome with agent incentives. Informally, a market outcome is *stable* if no pair of agents would rather match with each other than abide by the market outcome, and stable matchings can be computed when preferences are fully known.

In the context of large-scale matching platforms, however, the assumption that preferences are known breaks down. Platforms usually cannot have users report their complete preference profiles. Moreover, users may not even be aware of what their own preferences are. For example, a freelancer may not exactly know what types of projects they prefer until actually trying out specific ones. In reality, a data-driven platform is more likely to learn information about preferences from repeated feedback<sup>1</sup> over time. Two questions now emerge: In such marketplaces, how can stable matchings be learned? And what underlying structural assumptions are necessary for efficient learning to be possible?

To address these questions, we propose and investigate a model for learning stable matchings from noisy feedback. We model the platform's learning problem using stochastic multi-armed bandits, which lets us leverage the extensive body of work in the bandit literature to analyze the data efficiency of learning (see Lattimore and Szepesvári [30] for a textbook treatment). More specifically, our three main contributions are: (i) We develop an incentive-aware learning objective—Subset Instability—that captures the distance of a market outcome from equilibrium. (ii) Using Subset Instability as a measure of regret, we show that any “UCB-based” algorithm from the classical bandit literature can be adapted to this incentive-aware setting. (iii) We instantiate this idea for several families of preference structures to design efficient algorithms for incentive-aware learning. This helps elucidate how preference structure affects the complexity of learning stable matchings.

*Designing the learning objective.* Since mistakes are inevitable while exploring and learning, achieving exact stability at every timestep is an unattainable goal. To address this issue, we lean on approximation, focusing on learning market outcomes that are *approximately* stable. Thus, we need a metric that captures the distance of a market outcome from equilibrium.<sup>2</sup>

<sup>1</sup>Feedback might arise from explicit sources (e.g., riders rating drivers after a Lyft ride) or implicit sources (e.g., engagement metrics on an app); in either case, feedback is likely to be sparse and noisy.

<sup>2</sup>Previous work [17, 32] has investigated utility difference (i.e., the difference between the total utility achieved by the selected matching and the utility achieved by a stable matching) as a measure of regret. However, this does not capture distance from equilibrium in matching markets with monetary transfers (see Section 4) or without monetary transfers (see Section 6.3.1).

We introduce a notion for approximate stability that we call *Subset Instability*. Specifically, we define the Subset Instability of a market outcome to be the maximum difference, over all subsets  $\mathcal{S}$  of agents, between the total utility of the maximum weight matching on  $\mathcal{S}$  and the total utility of  $\mathcal{S}$  under the market outcome.<sup>3</sup> We show that Subset Instability can be interpreted as the amount the platform would have to *subsidize* participants to keep them on the platform and make the resulting matching stable. We can also interpret Subset Instability as the platform's cost of learning when facing competing platforms with greater knowledge of user preferences. Finally, we show that Subset Instability is the maximum gain in utility that a coalition of agents could have derived from an alternate matching such that no agent in the coalition is worse off.

Subset Instability also satisfies the following properties, which make it suitable for learning: (i) Subset Instability is equal to zero if and only if the market outcome is (exactly) stable; (ii) Subset Instability is robust to small perturbations to the utility functions of individual agents, which is essential for learning with noisy feedback; and (iii) Subset Instability upper bounds the utility difference of a market outcome from the socially optimal market outcome.

*Designing algorithms for learning a stable matching.* Using Subset Instability, we investigate the problem of learning a stable market outcome from noisy user feedback using the stochastic contextual bandit model (see, e.g., Lattimore and Szepesvári [30]). In each round, the platform selects a market outcome (i.e., a matching along with transfers) with the goal of minimizing cumulative instability.

We develop a general approach for designing bandit algorithms within our framework. Our approach is based on a primal-dual formulation of matching with transfers [42], in which the primal variables correspond to the matching and the dual variables can be used to set the transfers. We find that “optimism in the face of uncertainty,” the principle underlying many UCB-style bandit algorithms [5, 30], can be adapted to this primal-dual setting. The resulting algorithm is simple: Maintain upper confidence bounds on the agent utilities and compute, in each round, an optimal primal-dual pair in terms of these upper confidence bounds. The crux of the analysis is the following lemma, which bounds instability by the gap between the upper confidence bound and true utilities:

**LEMMA 1.1 (INFORMAL, SEE LEMMA 5.4 FOR A FORMAL STATEMENT).** *Given confidence sets for each utility value such that each confidence set contains the true utility, let  $(X, \tau)$  be a stable matching with transfers with respect to the utility functions given by the upper confidence bounds. The instability of  $(X, \tau)$  is upper bounded by the sum of the sizes of the confidence sets of pairs in  $X$ .*

We can thus analyze our algorithms by combining Lemma 1.1 with the analyses of existing UCB-style algorithms. In particular, we can essentially inherit the bounds on the size of the confidence bounds from traditional analyses of multi-arm bandits.

*Complexity of learning a stable matching.* Our main technical result is a collection of regret bounds for different structural assumptions on agent preferences. These bounds resemble the classical stochastic multi-armed bandits bounds when rewards have related structural assumptions. We summarize these regret bounds in Table 1 and elaborate on them in more detail below.

**THEOREM 1.2 (UNSTRUCTURED PREFERENCES, INFORMAL).** *For unstructured preferences, there exists a UCB-style algorithm that incurs  $\tilde{O}(N\sqrt{nT})$  regret according to Subset Instability after  $T$  rounds, where  $N$  is the number of agents on the platform and  $n$  is the number of agents that arrive in any round. (This bound is optimal up to logarithmic factors.)*

**THEOREM 1.3 (TYPED PREFERENCES, INFORMAL).** *Consider preferences such that each agent  $a$  has a type  $c_a \in \mathcal{C}$  and the utility of  $a$  when matched to another agent  $a'$  is given by a function of the*

<sup>3</sup>This formulation is inspired by the *strong  $\epsilon$ -core* of Shapley and Shubik [41].

Table 1. Regret Bounds for Varying Preference Structures (for  $N$  Agents and  $n$  Agents Per Round)

	Regret bound
Unstructured preferences	$\tilde{O}(N\sqrt{nT})$
Typed preferences	$\tilde{O}( C \sqrt{nT})$
Separable linear preferences	$\tilde{O}(d\sqrt{N}\sqrt{nT})$

types  $c_a$  and  $c_{a'}$ . There exists a UCB-style algorithm that incurs  $\tilde{O}(|C|\sqrt{nT})$  regret according to Subset Instability after  $T$  rounds, where  $n$  is the maximum number of agents that arrive to the platform in any round.

**THEOREM 1.4 (SEPARABLE LINEAR PREFERENCES, INFORMAL).** *Consider preferences such that the utility of an agent  $a$  when matched to another agent  $a'$  is  $\langle \phi(a), c_{a'} \rangle$ , where  $\phi(a) \in \mathbb{R}^d$  is unknown and  $c_{a'} \in \mathbb{R}^d$  is known. There exists a UCB-style algorithm that incurs  $\tilde{O}(d\sqrt{N}\sqrt{nT})$  regret according to Subset Instability after  $T$  rounds, where  $N$  is the number of agents on the platform and  $n$  is the maximum number of agents that arrive in any round.*

These results elucidate the role of preference structure on the complexity of learning a stable matching. Our regret bounds scale with  $N\sqrt{nT}$  for unstructured preferences (Theorem 1.2),  $|C|\sqrt{nT}$  for typed preferences (Theorem 1.3), and  $d\sqrt{N}\sqrt{nT}$  for linear preferences (Theorem 1.4). To illustrate these differences in a simple setting, let us consider the case where all of the agents show up every round, so  $n = N$ . In this case, our regret bound for unstructured preferences is superlinear in  $N$ ; in fact, this dependence on  $N$  is *necessary* as we demonstrate via a lower bound (see Lemma 5.5). However, the complexity of learning a stable matching changes substantially with preference structure assumptions. In particular, our regret bounds are sublinear/linear in  $N$  for typed preferences and separable linear preferences. This means that in large markets, a centralized platform can efficiently learn a stable matching with these preference structure assumptions.

*Connections and extensions.* Key to our results and extensions is the primal-dual characterization of equilibria in matching markets with transfers. Specifically, equilibria are described by a linear program whose primal form maximizes total utility over matchings and whose dual variables correspond to transfers. This linear program inspires our definition of Subset Instability, connects Subset Instability to platform profit (see Section 6.2), and relates learning with Subset Instability to regret minimization in combinatorial bandits (see Section 5.4). We adapt ideas from combinatorial bandits to additionally obtain  $O(\log T)$  instance-dependent regret bounds (see Section 6.1).

Our approach also offers a new perspective on learning stable matchings in markets with *non-transferable* utilities [17, 32]. Although this setting does not admit a linear program formulation, we show Subset Instability can be extended to what we call NTU Subset Instability (see Section 6.3), which turns out to have several advantages over the instability measures studied in previous work. Our algorithmic principles extend to NTU Subset Instability: We prove regret bounds commensurate with those for markets with transferable utilities.

## 1.1 Related Work

In the machine learning literature, starting with Das and Kamenica [17] and Liu et al. [32], several works [9, 12, 17, 32, 33, 40] study learning stable matchings from bandit feedback in the Gale-Shapley stable marriage model [23]. A major difference between this setting and ours is the absence of monetary transfers between agents. These works focus on the *utility difference* rather

than the instability measure that we consider. Cen and Shah [12] extend this bandits model to incorporate fixed, predetermined cost/transfer rules. However, they do not allow the platform to set arbitrary transfers between agents. Moreover, they also consider a weaker notion of stability that does not consider agents negotiating arbitrary transfers: Defecting agents must set their transfers according to a fixed, predetermined structure. In contrast, we follow the classical definition of stability [42].

Outside of the machine learning literature, several papers also consider the complexity of finding stable matchings in other feedback and cost models, e.g., communication complexity [4, 24, 43] and query complexity [4, 20]. Of these works, Shi [43], which studies the communication complexity of finding approximately stable matchings with transferable utilities, is perhaps most similar to ours. This work assumes agents know their preferences and focuses on the communication bottleneck, whereas we study the costs associated with learning preferences. Moreover, the approximate stability notion in Shi [43] is the maximum unhappiness of any *pair* of agents, whereas Subset Instability is equivalent to the maximum unhappiness over any *subset* of agents. For learning stable matchings, Subset Instability has the advantages of being more fine-grained and having a primal view that motivates a clean UCB-based algorithm.

Our notion of instability connects to historical works in coalitional game theory: Related are the concepts of the strong- $\epsilon$  core of Shapley and Shubik [41] and the indirect function of Martínez-Legaz [37], although each was introduced in a very different context than ours. Nonetheless, they reinforce the fact that our instability notion is a very natural one to consider.

A complementary line of work in economics [2, 11, 34, 35] considers stable matchings under incomplete information. These works focus on defining stability when the agents have incomplete information about their own preferences, whereas we focus on the platform's problem of learning stable matchings from noisy feedback. As a result, these works relax the definition of stability to account for uncertainty in the preferences of agents, rather than the uncertainty experienced by the platform from noisy feedback.

Multi-armed bandits have also been applied to learning in other economic contexts. For example, learning a socially optimal matching (without learning transfers) is a standard application of combinatorial bandits [13, 15, 16, 22, 29]. Other applications at the interface of bandit methodology and economics include dynamic pricing [8, 27, 38], incentivizing exploration [21, 36], learning under competition [3], and learning in matching markets without incentives [26].

Finally, primal-dual methods have also been applied to other problems in the bandits literature (e.g., References [25, 31, 44]).

## 2 PRELIMINARIES

The foundation of our framework is the *matching with transfers* model of Shapley and Shubik [42]. In this section, we introduce this model along with the concept of stable matching.

### 2.1 Matching with Transferable Utilities

Consider a two-sided market that consists of a finite set  $\mathcal{I}$  of customers on one side and a finite set  $\mathcal{J}$  of providers on the other. Let  $\mathcal{A} := \mathcal{I} \cup \mathcal{J}$  be the set of all agents. A *matching*  $X \subseteq \mathcal{I} \times \mathcal{J}$  is a set of pairs  $(i, j)$  that are pairwise disjoint, representing the pairs of agents that are matched. Let  $\mathcal{X}_{\mathcal{A}}$  denote the set of all matchings on  $\mathcal{A}$ . For notational convenience, we define for each matching  $X \in \mathcal{X}_{\mathcal{A}}$  an equivalent functional representation  $\mu_X: \mathcal{A} \rightarrow \mathcal{A}$ , where  $\mu_X(i) = j$  and  $\mu_X(j) = i$  for all matched pairs  $(i, j) \in X$ , and  $\mu_X(a) = a$  if  $a \in \mathcal{A}$  is unmatched.

When a pair of agents  $(i, j) \in \mathcal{I} \times \mathcal{J}$  matches, each experiences a utility gain. We denote these utilities by a global utility function  $u: \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}$ , where  $u(a, a')$  denotes the utility that agent  $a$  gains from being matched to agent  $a'$ . (If  $a$  and  $a'$  are on the same side of the market, then we



take  $u(a, a')$  to be zero by default.) We allow these utilities to be negative if matching results in a net cost (e.g., if an agent is providing a service). We assume each agent  $a \in \mathcal{A}$  receives zero utility if unmatched, i.e.,  $u(a, a) = 0$ . When we wish to emphasize the role of an individual agent's utility function, we will use the equivalent notation  $u_a(a') := u(a, a')$ .

A *market outcome* consists of a matching  $X \in \mathcal{X}_{\mathcal{A}}$  along with a vector  $\tau \in \mathbb{R}^{\mathcal{A}}$  of transfers, where  $\tau_a$  is the amount of money transferred from the platform to agent  $a$  for each  $a \in \mathcal{A}$ . These monetary transfers are a salient feature of most real-world matching markets: Riders pay drivers on Lyft, clients pay freelancers on TaskRabbit, and guests pay hosts on Airbnb. Shapley and Shubik [42] capture this aspect of matching markets by augmenting the classical two-sided matching model with transfers of utility between agents. Transfers are typically required to be *zero-sum*, meaning that  $\tau_i + \tau_j = 0$  for all matched pairs  $(i, j) \in X$  and  $\tau_a = 0$  if  $a$  is unmatched. Here,  $X$  represents how agents are matched and  $\tau_a$  represents the transfer that agent  $a$  receives (or pays). The net utility that an agent  $a$  derives from a matching with transfers  $(X, \tau)$  is therefore  $u(a, \mu_X(a)) + \tau_a$ .

*Stable matchings.* In matching theory, stability captures when a market outcome aligns with individual agents' preferences. Roughly speaking, a market outcome  $(X, \tau)$  is stable if: (i) no individual agent  $a$  would rather be unmatched and (ii) no pair of agents  $(i, j)$  can agree on a transfer such that both would rather match with each other than abide by  $(X, \tau)$ . Formally:

*Definition 2.1.* A market outcome  $(X, \tau)$  is *stable* if: (i) it is *individually rational*, i.e.,

$$u_a(\mu_X(a)) + \tau_a \geq 0 \quad (1)$$

for all agents  $a \in \mathcal{A}$ , and (ii) it *has no blocking pairs*, i.e.,

$$(u_i(\mu_X(i)) + \tau_i) + (u_j(\mu_X(j)) + \tau_j) \geq u_i(j) + u_j(i) \quad (2)$$

for all pairs of agents  $(i, j) \in \mathcal{I} \times \mathcal{J}$ .<sup>4</sup>

A fundamental property of the matching with transfers model is that if  $(X, \tau)$  is stable, then  $X$  is a maximum weight matching, i.e.,  $X$  maximizes  $\sum_{a \in \mathcal{A}} u_a(\mu_X(a))$  over all matchings  $X \in \mathcal{X}_{\mathcal{A}}$  [42]. The same work shows that stable market outcomes coincide with Walrasian equilibria. (For completeness, we recapitulate the basic properties of this model in Appendix A.)

To make the matching with transfers model concrete, we use the simple market depicted in the center panel of Figure 1 as a running example throughout the article. This market consists of a customer Charlene and two providers Percy and Quinn, which we denote by  $\mathcal{I} = \{C\}$  and  $\mathcal{J} = \{P, Q\}$ . If the agents' utilities are as given in Figure 1, then Charlene would prefer Quinn, but Quinn's cost of providing the service is much higher. Thus, matching Charlene and Percy is necessary for a stable outcome. This matching is stable for any transfer from Charlene to Percy in the interval [5, 7].

### 3 LEARNING PROBLEM AND FEEDBACK MODEL

We instantiate the platform's learning problem in a stochastic contextual bandits framework. Matching takes place over the course of  $T$  rounds. We denote the set of all customers by  $\mathcal{I}^*$ , the set of all providers by  $\mathcal{J}^*$ , and the set of all agents on the platform by  $\mathcal{A}^* = \mathcal{I}^* \cup \mathcal{J}^*$ . Each agent  $a \in \mathcal{A}^*$  has an associated context  $c_a \in \mathcal{C}$ , where  $\mathcal{C}$  is the set of all possible contexts. This context represents the side information available to the platform about the agent, e.g., demographic,

<sup>4</sup>We observe that (2) corresponds to no pair of agents  $(i, j)$  being able to agree on a transfer such that both would rather match with each other than abide by  $(X, \tau)$ . Notice that a pair  $(i, j)$  violates (2) if and only if they can find a transfer  $\tau'_i = -\tau'_j$  such that  $u_i(j) + \tau'_i > u_i(\mu_X(i)) + \tau_i$  and  $u_j(i) + \tau'_j > u_j(\mu_X(j)) + \tau_j$ .

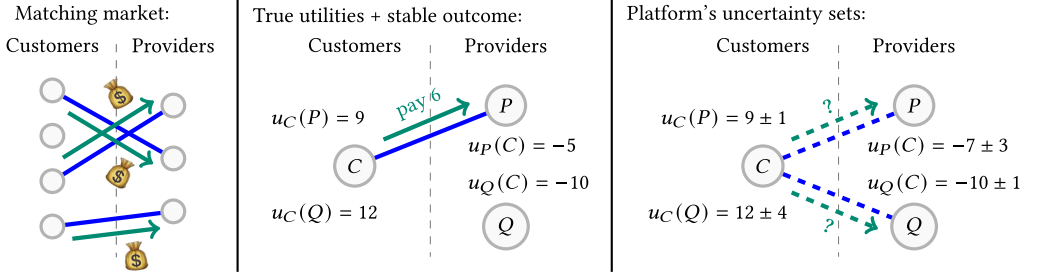


Fig. 1. The left panel depicts a schematic of a matching (blue) with transfers (green). The center panel depicts a matching market with three agents and a stable matching with transfers for that market. (If the transfer 6 is replaced with any value between 5 and 7, then the outcome remains stable.) The right panel depicts the same market, but with utilities replaced by uncertainty sets; note that no matching with transfers is stable for all realizations of utilities.

location, or platform usage information. Each round, a set of agents arrives to each side of the market. The platform then selects a market outcome and incurs a regret equal to the *instability* of the market outcome (which we introduce formally in Section 4). Finally, the platform receives noisy feedback about the utilities of each matched pair  $(i, j)$ .

To interpret the noisy feedback, note that platforms in practice often receive feedback both explicitly (e.g., riders rating drivers after a Lyft ride) and implicitly (e.g., engagement metrics on an app). In either instance, feedback is likely to be sparse and noisy. For simplicity, we do not account for agents strategically manipulating their feedback to the platform and focus on the problem of learning preferences from unbiased reports.

We now describe this model more formally. In the  $t$ -th round:

- (1) A set  $I^t \subseteq I^*$  of customers and a set  $J^t \subseteq J^*$  of providers arrive to the market. Write  $I^t \cup J^t =: \mathcal{A}^t$ . The platform observes the identity  $a$  and the *context*  $c_a \in C$  of each agent  $a \in \mathcal{A}^t$ .
- (2) The platform selects a matching with *zero-sum* transfers  $(X^t, \tau^t)$  between  $I^t$  and  $J^t$ .
- (3) The platform observes noisy utilities  $u_a(\mu_{X^t}(a)) + \varepsilon_{a,t}$  for each agent  $a \in I^t \cup J^t$ , where the  $\varepsilon_{a,t}$  are independent, 1-subgaussian random variables.<sup>5</sup>
- (4) The platform incurs regret equal to the *instability* of the selected market outcome  $(X^t, \tau^t)$ . (We define instability formally in Section 4.)

The platform's total regret  $R_T$  is thus the cumulative instability incurred up through round  $T$ .

### 3.1 Preference Structure

In this bandits framework, we can impose varying degrees of structure on agent preferences. We encode these preference structures via the functional form of agents' utility functions and their relation to agent contexts. More formally, let  $\mathcal{U}$  be the set of functions  $u: \mathcal{A}^* \times \mathcal{A}^* \rightarrow \mathbb{R}$ , i.e.,  $\mathcal{U}$  is the set of all possible (global) utility functions. We now introduce several classes of preference structures as subsets of  $\mathcal{U}$ .

*Unstructured preferences.* The simplest setting we consider is one where the preferences are unstructured. Specifically, we consider the class of utility functions

$$\mathcal{U}_{\text{unstructured}} = \{u \in \mathcal{U} \mid u(a, a') \in [-1, 1]\}.$$

<sup>5</sup>Our feedback model corresponds to *semi-bandit* feedback, since the platform has (noisy) access to each agent's utility within the matching rather than the overall utility of the matching.

(Here, one can think of the context as being uninformative, i.e.,  $C$  is the singleton set.) In this setup, the platform must learn each agent's utility function  $u_a(\cdot) = u(a, \cdot)$ .

*Typed preferences.* We next consider a market where each agent comes in one of finitely many types, with agents of the same type having identical preferences. Assuming typed preference structures is standard in theoretical models of markets (see, e.g., Azevedo and Hatfield [7], Debreu and Scarf [18], Echenique et al. [19]). We can embed types into our framework by having each agent's context represent their type, with  $|C| < \infty$ . The global utility function is then fully specified by agents' contexts:

$$\mathcal{U}_{\text{typed}} = \{u \in \mathcal{U} \mid u(a, a') = f(c_a, c_{a'}) \text{ for some } f: C \times C \rightarrow [-1, 1]\}.$$

*Separable linear preferences.* We next consider markets where each agent is associated with *known* information given by their context as well as *hidden* information that must be learned by the platform. (This differs from unstructured preferences, where all information was hidden, and typed preferences, where each agent's context encapsulated their full preferences.) We explore this setting under the assumption that agents' contexts and hidden information interact linearly.

We assume that all contexts belong to  $\mathcal{B}^d$  (i.e.,  $C = \mathcal{B}^d$ ) where  $\mathcal{B}^d$  is the  $\ell_2$  unit ball in  $\mathbb{R}^d$ . We also assume that there exists a function  $\phi: \mathcal{A}^* \rightarrow \mathcal{B}^d$  mapping each agent to the hidden information associated to that agent. The preference class  $\mathcal{U}_{\text{linear}}^d$  can then be defined as

$$\mathcal{U}_{\text{linear}}^d = \{u \in \mathcal{U} \mid u(a, a') = \langle c_{a'}, \phi(a) \rangle \text{ for some } \phi: \mathcal{A}^* \rightarrow \mathcal{B}^d\}.$$

#### 4 MEASURING APPROXIMATE STABILITY

When learning stable matchings, we must settle for guarantees of approximate stability, since exact stability—a binary notion—is unattainable when preferences are uncertain. To see this, we return to the example from Figure 1. Suppose that the platform has uncertainty sets given by the right panel. Recall that for the true utilities, all stable outcomes match Charlene with Percy. If the true utilities were instead the upper bounds of each uncertainty set, then all stable outcomes would match Charlene and Quinn. Given only the uncertainty sets, it is impossible for the platform to find an (exactly) stable matching, so it is necessary to introduce a measure of approximate stability as a relaxed benchmark for the platform; we turn to this now.

Given the insights of Shapley and Shubik [42]—that all stable outcomes maximize the sum of agents' utilities—it might seem natural to measure distance from stability simply in terms of the *utility difference*. To define this formally, let  $\mathcal{A}$  be the set of agents participating in the market. (This corresponds to  $\mathcal{A}^t$  at timestep  $t$  in the bandits model.) The utility difference<sup>6</sup> of a market outcome  $(X, \tau)$  is given by:

$$\left( \max_{X' \in \mathcal{X}_{\mathcal{A}}} \sum_{a \in \mathcal{A}} u_a(\mu_{X'}(a)) \right) - \left( \sum_{a \in \mathcal{A}} u_a(\mu_X(a)) + \tau_a \right). \quad (3)$$

The first term  $\max_{X' \in \mathcal{X}_{\mathcal{A}}} \sum_{a \in \mathcal{A}} u_a(\mu_{X'}(a))$  is the maximum total utility of any matching, and the second term  $\sum_{a \in \mathcal{A}} (u_a(\mu_X(a)) + \tau_a)$  is the total utility of market outcome  $(X, \tau)$ . Since transfers are zero-sum, (3) can be equivalently written as

$$\left( \max_{X' \in \mathcal{X}_{\mathcal{A}}} \sum_{a \in \mathcal{A}} u_a(\mu_{X'}(a)) \right) - \sum_{a \in \mathcal{A}} u_a(\mu_X(a)).$$

<sup>6</sup>Utility difference is standard as a measure of regret for learning a maximum weight matching in the combinatorial bandits literature (see, e.g., Gai et al. [22]). However, we show that for learning stable matchings, a fundamentally different measure of regret is needed.



But this shows that utility difference actually ignores the transfers  $\tau$  entirely! In fact, the utility difference can be zero even when the transfers lead to a market outcome that is far from stable (see Appendix B.1). Utility difference is therefore *not* incentive-aware, making it unsuitable as an objective for learning stable matchings with transfers.

In the remainder of this section, we propose a measure of instability—Subset Instability—which we will show serves as a suitable objective for learning stable matchings with transfers. Specifically, we show that Subset Instability captures the distance of a market outcome from equilibrium while reflecting both the platform’s objective and the users’ incentives. We additionally show that Subset Instability satisfies several structural properties that make it useful for learning.

#### 4.1 Subset Instability

Subset Instability is based on utility difference, but rather than only looking at the market in aggregate, it takes a maximum ranging over all subsets of agents.

*Definition 4.1.* Given utilities  $u$ , the *Subset Instability*  $I(X, \tau; u, \mathcal{A})$  of a matching with transfers  $(X, \tau)$  is

$$\max_{S \subseteq \mathcal{A}} \left[ \left( \max_{X' \in \mathcal{X}_S} \sum_{a \in S} u_a(\mu_{X'}(a)) \right) - \left( \sum_{a \in S} u_a(\mu_X(a)) + \tau_a \right) \right]. \quad (*)$$

(The first term  $\max_{X' \in \mathcal{X}_S} \sum_{a \in S} u_a(\mu_{X'}(a))$  is the maximum total utility of any matching over  $S$ , and the second term  $\sum_{a \in S} (u_a(\mu_X(a)) + \tau_a)$  is the total utility of the agents in  $S$  under market outcome  $(X, \tau)$ .)

Intuitively, Subset Instability captures stability because it checks whether any subset of agents would prefer an alternate outcome. We provide a more extensive economic interpretation below; but before doing so, we first illustrate Definition 4.1 in the context of the example in Figure 1.

Consider the matching  $X = \{(C, Q)\}$  with transfers  $\tau_C = -11$  and  $\tau_Q = 11$ . (This market outcome is stable for the upper bounds of the uncertainty sets of the platform in Figure 1, but not stable for the true utilities.) It is not hard to see that the subset  $S$  that maximizes Subset Instability is  $S = \{C, P\}$ , in which case  $\max_{X' \in \mathcal{X}_S} \sum_{a \in S} u_a(\mu_{X'}(a)) = 4$  and  $\sum_{a \in S} (u_a(\mu_X(a)) + \tau_a) = 1$ . Thus, the Subset Instability of  $(X, \tau)$  is  $I(X, \tau; u, \mathcal{A}) = 4 - 1 = 3$ . In contrast, the utility difference of  $(X, \tau)$  is 2.

We now discuss several interpretations of Subset Instability, which provide further insight into why Subset Instability serves as a meaningful notion of approximate stability in online marketplaces. In particular, Subset Instability can be interpreted as *the minimum stabilizing subsidy*, as *the platform’s cost of learning*, as *a measure of user unhappiness*, and as *a distance from equilibrium*.

*Subset Instability as the platform’s minimum stabilizing subsidy.* Subset Instability can be interpreted in terms of monetary subsidies from the platform to the agents. Specifically, the Subset Instability of a market outcome equals the minimum amount the platform could subsidize agents so the subsidized market outcome is individually rational and has no blocking pairs.

More formally, let  $s \in \mathbb{R}_{\geq 0}^{\mathcal{A}}$  denote subsidies made by the platform, where the variable  $s_a \geq 0$  represents the subsidy provided to agent  $a$ .<sup>7</sup> For a market outcome  $(X, \tau)$ , the *minimum stabilizing subsidy* is

$$\min_{s \in \mathbb{R}_{\geq 0}^{\mathcal{A}}} \left\{ \sum_{a \in \mathcal{A}} s_a \mid (X, \tau + s) \text{ is stable} \right\}, \quad (4)$$

<sup>7</sup>The requirement that  $s_a \geq 0$  enforces that all subsidies are nonnegative; without it, (5) would reduce to the utility difference, which is not incentive-aware.

where we define stability in analogy to Definition 2.1. Specifically, we say that a market outcome  $(X, \tau)$  with subsidies  $s$  is *stable* if it is individually rational, i.e.,  $u_a(\mu_X(a)) + \tau_a + s_a \geq 0$  for all agents  $a \in \mathcal{A}$ , and has no blocking pairs, i.e.,  $(u_i(\mu_X(i)) + \tau_i + s_i) + (u_j(\mu_X(j)) + \tau_j + s_j) \geq u_i(j) + u_j(i)$  for all pairs of agents  $(i, j) \in \mathcal{I} \times \mathcal{J}$ .

Given this setup, we show the following equivalence:

PROPOSITION 4.2. *Minimum stabilizing subsidy equals Subset Instability for any market outcome.*

The proof boils down to showing that the two definitions are “dual” to each other. To formalize this, we rewrite the minimum stabilizing subsidy as the solution to the following linear program:<sup>8</sup>

$$\begin{aligned} \min_{s \in \mathbb{R}^{|\mathcal{A}|}} \quad & \sum_{a \in \mathcal{A}} s_a \\ \text{s.t.} \quad & (u_i(\mu_X(i)) + \tau_i + s_i) + (u_j(\mu_X(j)) + \tau_j + s_j) \geq u_i(j) + u_j(i) & \forall (i, j) \in \mathcal{I} \times \mathcal{J} \\ & u_a(\mu_X(a)) + \tau_a + s_a \geq 0 & \forall a \in \mathcal{A} \\ & s_a \geq 0 & \forall a \in \mathcal{A}. \end{aligned} \tag{5}$$

The crux of our argument is that the dual linear program to (5) maximizes the combinatorial objective (\*). The equivalence of (\*) and (5) then follows from strong duality.

With this alternate formulation of Subset Instability in mind, we revisit the example in Figure 1. Again, consider the matching  $X = \{(C, Q)\}$  with transfers  $\tau_C = -11$  and  $\tau_Q = 11$ . (This is stable for the upper bounds of the uncertainty sets of the platform in Figure 1, but not stable for the true utilities.) We have already shown above that the Subset Instability of this market outcome is 3. To see this via the subsidy formulation, note that the optimal subsidy  $s$  gives  $C$  and  $P$  a total of 3. (E.g., we give  $C$  a subsidy of  $s_C = 2$  and  $P$  a subsidy of  $s_P = 1$ .) Indeed, if  $s_C + s_P = 3$ , then

$$(u_C(\mu_X(C)) + \tau_C + s_C) + (u_P(\mu_X(P)) + \tau_P + s_P) \geq u_C(P) + u_P(C)$$

holds (with equality), so the pair  $(C, P)$  could no longer gain by matching with each other.

The subsidy perspective turns out to be useful when designing learning algorithms. In particular, while the formulation in Definition 4.1 involves a maximization over the  $2^{|\mathcal{A}|}$  subsets of  $\mathcal{A}$ , the linear programming formulation (5) only involves  $O(|\mathcal{A}|)$  variables and  $O(|\mathcal{A}|^2)$  constraints.

*Subset Instability as the platform’s cost of learning.* We next connect minimum stabilizing subsidies to the platform’s *cost of learning*—how much the platform would have to pay to keep users on the platform in the presence of a worst-case (but budget-balanced) competitor with perfect knowledge of agent utilities.

Observe that (4) is the minimum amount the platform could subsidize agents so that no budget-balanced competitor could convince agents to leave. The way that we formalize “convincing agents to leave” is that: (a) an agent will leave the original platform if they prefer to be unmatched over being on the platform or (b) a pair of agents who are matched on the competitor’s platform will leave the original platform if they both prefer the new market outcome over their original market outcomes. Thus, if we imagine the platform as actually paying the subsidies, then the cumulative instability (i.e., our regret) can be realized as a “cost of learning”: It is how much the platform pays the agents to learn a stable outcome while ensuring that no agent has the incentive to leave during the learning process. Later on, we will see that our algorithmic approach can be extended to efficiently compute feasible subsidies for (5) that are within a constant factor of our regret bound, meaning that subsidies can be implemented using only the information that the platform

<sup>8</sup>In this linear program, the first set of constraints ensures there are no blocking pairs, while the second set of constraints ensures individual rationality.

has. Moreover, in Section 6.2, we show that cost of learning can also be explicitly connected to the platform's revenue.

*Subset Instability as a measure of user unhappiness.* While the above interpretations focus on Subset Instability from the platform's perspective, we show that Subset Instability can also be interpreted as a measure of user unhappiness. Given a subset  $\mathcal{S} \subseteq \mathcal{A}$  of agents, which we call a coalition, we define the *unhappiness* of  $\mathcal{S}$  with respect to a market outcome  $(X, \tau)$  to be the maximum gain (relative to  $(X, \tau)$ ) in total utility that the members of coalition  $\mathcal{S}$  could achieve by matching only among themselves, such that no member is worse off than they were in  $(X, \tau)$ . (See Appendix B.3 for a formal definition.) The condition that no member is worse off ensures that all agents would actually want to participate in the coalition (i.e. they prefer it to the original market outcome).

User unhappiness differs from the original definition of Subset Instability in (\*), because (\*) does not require individuals to be better off in any alternative matching. However, we show that this difference is inconsequential:

**PROPOSITION 4.3.** *The maximum unhappiness of any coalition  $\mathcal{S} \subseteq \mathcal{A}$  with respect to  $(X, \tau)$  equals the Subset Instability  $I(X, \tau; u, \mathcal{A})$ .*

See Appendix B.3 for a full proof. In the proof, we relate the maximum unhappiness of any coalition to the dual linear program to (5). To show this relation, we leverage the fact that optimal solutions to the dual program correspond to blocking pairs of agents as well as individual rationality violations.

The main takeaway from Proposition 4.3 is that Subset Instability not only measures costs to the platform, but also costs to users, in terms of the maximum amount they “leave on the table” by not negotiating an alternate arrangement amongst themselves.

*Subset Instability as a distance from equilibrium.* Finally, we connect Subset Instability to solution concepts for coalitional games, a general concept in game theory that includes matching with transfers as a special case. Coalitional games (also known as cooperative games) capture competition and cooperation amongst a group of agents. The *core* is the set of outcomes in a cooperative game such that no subset  $\mathcal{S}$  of agents can achieve higher total utility among themselves than according to the given outcome. In games where the core is empty, a natural relaxation is the *strong  $\varepsilon$ -core* [41], which is the set of outcomes in a cooperative game such that no subset  $\mathcal{S}$  of agents can achieve total utility among themselves that is at least  $\varepsilon$  greater than according to the given outcome.

Subset Instability can be seen as transporting the strong  $\varepsilon$ -core notion to a slightly different context. In particular, in the context of matching with transferable utilities, the core is exactly the set of stable matchings; since a stable matching always exists, the core is always nonempty. Even though the core is nonempty, we can nonetheless use the strong  $\varepsilon$ -core to measure *distance from the core*. More specifically, it is natural to consider the smallest  $\varepsilon$  such that  $(X, \tau)$  is in the strong  $\varepsilon$ -core. This definition exactly aligns with Subset Instability, thus providing an alternate interpretation of Subset Instability within the context of coalitional game theory.

## 4.2 Properties of Subset Instability

We now describe additional properties of our instability measure that are important for learning. We show that Subset Instability is: (i) zero if and only if the matching with transfers is stable, (ii) Lipschitz in the true utility functions, and (iii) lower bounded by the utility difference.

PROPOSITION 4.4. *Subset Instability satisfies the following properties:*

- (1) *Subset Instability is always nonnegative and is zero if and only if  $(X, \tau)$  is stable.*
- (2) *Subset Instability is Lipschitz continuous with respect to agent utilities. That is, for any possible market outcome  $(X, \tau)$  and any pair of utility functions  $u$  and  $\tilde{u}$  it holds that:*

$$|I(X, \tau; u, \mathcal{A}) - I(X, \tau; \tilde{u}, \mathcal{A})| \leq 2 \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_{\infty}.$$

- (3) *Subset Instability is always at least the utility difference.*

We defer the proof to Appendix B.4.

These three properties show that Subset Instability is useful as a regret measure for learning stable matchings. The first property establishes that Subset Instability satisfies the basic desideratum of having zero instability coincide with exact stability. The second property shows that Subset Instability is robust to small perturbations to the utility functions of individual agents. The third property ensures that, when learning using Subset Instability as a loss function, the platform learns a socially optimal matching.

Note that the second property already implies the existence of an explore-then-commit algorithm that achieves  $\tilde{O}(N^{4/3}T^{2/3})$  regret in the simple setting where  $\mathcal{A}^t = \mathcal{A}$  for some  $\mathcal{A}$  of size  $N$  for all  $t$ .<sup>9</sup> In the next section, we will explore algorithms that improve the dependence on the number of rounds  $T$  to  $\sqrt{T}$  and also work in more general settings.

## 5 REGRET BOUNDS

In this section, we develop a general approach for designing algorithms that achieve near-optimal regret within our framework. To be precise, the platform's regret is defined to be

$$R_T = \sum_{t=1}^T I(X^t, \tau^t; u, \mathcal{A}^t).$$

While our framework bears some resemblance to the (incentive-free) combinatorial bandit problem of learning a maximum weight matching, two crucial differences differentiate our setting: (i) In each round, the platform must choose *transfers* in addition to a matching, and (ii) loss is measured with respect to *instability* rather than the utility difference. Nonetheless, we show that a suitable interpretation of “optimism in the face of uncertainty” can still apply.

*Regret bounds for different preference structures.* By instantiating this optimism-based approach, we derive regret bounds for the preference structures introduced in Section 3. We start with the simplest case of unstructured preferences, where we assume no structure on the utilities.

**THEOREM 5.1.** *For preference class  $\mathcal{U}_{\text{unstructured}}$  (see Section 3), MATCHUCB (defined in Section 5.3) incurs expected regret  $\mathbb{E}(R_T) = O(|\mathcal{A}|\sqrt{nT} \log(|\mathcal{A}|T))$ , where  $n = \max_t |\mathcal{A}_t|$ .*

In Section 5.4, we additionally give a matching (up to logarithmic factors) lower bound showing for  $n = |\mathcal{A}|$  that such scaling in  $|\mathcal{A}|$  is indeed necessary. This demonstrates that the regret scales with  $|\mathcal{A}|\sqrt{n}$ , which is superlinear in the size of the market. Roughly speaking, this bound means that the platform is required to learn a superconstant amount of information per agent in the

<sup>9</sup>This bound can be achieved by adapting the **explore-then-commit (ETC)** approach, where the platform explores by choosing each pair of agents  $\tilde{O}((T/N)^{2/3})$  times [30]. Thus,  $\tilde{O}(N^{1/3}T^{2/3})$  rounds are spent exploring, and the Subset Instability of the matching selected in the commit phase is  $\tilde{O}(N^{4/3}T^{2/3})$  with high probability. We omit further details, since this analysis is a straightforward adaptation of the typical ETC analysis.

marketplace. These results suggest that without preference structure, it is unlikely that a platform can efficiently learn a stable matching in large markets.

The next two bounds demonstrate that, with preference structure, efficient learning of a stable matching becomes possible. First, we consider typed preferences, which are purely specified by a function  $f$  mapping finitely many pairs of contexts to utilities.

**THEOREM 5.2.** *For preference class  $\mathcal{U}_{\text{typed}}$  (see Section 3), MATCHTYPEDUCB (defined in Section 5.3) incurs expected regret  $\mathbb{E}(R_T) = O(|C|\sqrt{nT \log(|\mathcal{A}|T)})$ , where  $n = \max_t |\mathcal{A}_t|$ .*

For a fixed type space  $C$ , the regret bound in Theorem 5.2 scales sublinearly with the market size (captured by  $|\mathcal{A}|$  and  $n$ ). This demonstrates that the platform can efficiently learn a stable matching when preferences are determined by types. In fact, the regret bound only depends on the number of agents who arrive on the platform in any round; notably, it does not depend on the total number of agents on the platform (beyond logarithmic factors).

Finally, we consider separable linear preferences, where the platform needs to learn hidden information associated with each agent.

**THEOREM 5.3.** *For preference class  $\mathcal{U}_{\text{linear}}$  (see Section 3), MATCHLINUCB (defined in Section 5.3) incurs expected regret  $\mathbb{E}(R_T) = O(d\sqrt{|\mathcal{A}|\sqrt{nT \log(|\mathcal{A}|T)})}$ , where  $n = \max_t |\mathcal{A}_t|$ .*

When  $n$  is comparable to  $|\mathcal{A}|$ , the regret bound in Theorem 5.3 scales linearly with the market size (captured by  $|\mathcal{A}|$ ) and linearly with the dimension  $d$ . Roughly speaking, this means that the platform learns (at most) a constant amount of information per agent in the marketplace. We interpret this as indicating that the platform can efficiently learn a stable matching in large markets for separable linear preferences, although learning in this setting is more demanding than for typed preferences.

## 5.1 Algorithm

Following the principle of optimism, our algorithm selects at each round a stable market outcome using upper confidence bounds as if they were the true agent utilities. To design and analyze this algorithm, we leverage the fact that, in the full-information setting, stable market outcomes are optimal solutions to a pair of primal-dual linear programs whose coefficients depend on agents' utility functions. This primal-dual perspective lets us compute a market outcome each round. A particular consequence is that any UCB-based algorithm for learning matchings in a semi-bandit setting can be transformed into an algorithm for learning *both* the matching and the prices.

*Stable market outcomes via linear programming duality.* Before proceeding with the details of our algorithm, we review how the primal-dual framework can be used to select a stable market outcome in the full information setting. Shapley and Shubik [42] show that stable market outcomes  $(X, \tau)$  correspond to optimal primal-dual solutions to the following pair of primal and dual linear programs (where we omit the round index  $t$  and consider matchings over  $\mathcal{A} = \mathcal{I} \cup \mathcal{J}$ ):

Primal (P)	Dual (D)
$\max_{Z \in \mathbb{R}^{ \mathcal{I}  \times  \mathcal{J} }} \sum_{(i,j) \in \mathcal{I} \times \mathcal{J}} Z_{i,j} (u_i(j) + u_j(i))$	$\min_{p \in \mathbb{R}^{ \mathcal{A} }} \sum_{a \in \mathcal{A}} p_a$
$\text{s.t. } \sum_{j \in \mathcal{J}} Z_{i,j} \leq 1 \quad \forall i \in \mathcal{I}$	$\text{s.t. } p_i + p_j \geq u_i(j) + u_j(i) \quad \forall (i,j) \in \mathcal{I} \times \mathcal{J}$
$\sum_{i \in \mathcal{I}} Z_{i,j} \leq 1 \quad \forall j \in \mathcal{J}$	$p_a \geq 0 \quad \forall a \in \mathcal{A}$
$Z_{i,j} \geq 0 \quad \forall (i,j) \in \mathcal{I} \times \mathcal{J}$	

**ALGORITHM 1:** METAMATCHUCB: A bandit meta-algorithm for matching with transferable utilities.

---

```

1: procedure METAMATCHUCB( $T$ )
2:   Initialize confidence intervals  $\mathcal{C}$  over utilities.
3:   for  $1 \leq t \leq T$  do
4:      $(X^t, \tau^t) \leftarrow \text{COMPUTEMATCH}(\mathcal{C})$ 
5:     Update confidence intervals  $\mathcal{C}$ .

```

---

**ALGORITHM 2:** COMPUTEMATCH: Compute matching with transfers from confidence sets of utilities

---

```

1: procedure COMPUTEMATCH( $\mathcal{C}$ )
2:   for  $(i, j) \in \mathcal{I} \times \mathcal{J}$  do ▷ Instantiate UCB estimates of utilities.
3:      $u_i^{\text{UCB}}(j) \leftarrow \max(C_{i,j})$ 
4:      $u_j^{\text{UCB}}(i) \leftarrow \max(C_{j,i})$ 
5:    $(X^*, p^*) \leftarrow$  optimal primal-dual pair for (P) and (D) given utilities  $u^{\text{UCB}}$ 
6:   for  $a \in \mathcal{A}$  do ▷ Set transfers based on  $(X^*, p^*)$  and UCB utilities.
7:      $\tau_a \leftarrow p_a^* - u_a^{\text{UCB}}(\mu_{X^*}(a))$ 
8:   return  $(X^*, \tau)$ 

```

---

The primal program (P) is a linear programming formulation of the maximum weight matching problem: The Birkhoff-von Neumann theorem states that its extreme points are exactly the indicator vectors for matchings between  $\mathcal{I}$  and  $\mathcal{J}$ . Each dual variable  $p_a$  in (D) can be interpreted as a *price* that roughly corresponds to agent  $a$ 's net utility. Specifically, given any optimal primal-dual pair  $(Z, p)$ , one can recover a matching  $\mu_X$  from the nonzero entries of  $Z$  and set transfers  $\tau_a = p_a - u_a(\mu_X(a))$  to obtain a stable outcome  $(X, \tau)$ . Moreover, any stable outcome induces an optimal primal-dual pair  $(Z, p)$ .

*Algorithm overview.* Leveraging the above primal-dual formulation of stability, we introduce a meta-algorithm METAMATCHUCB for learning stable outcomes (Algorithm 1). In each round, we compute a matching with transfers by solving the primal-dual linear programs for our upper confidence bounds: Suppose we have a collection  $\mathcal{C}$  of confidence sets  $C_{i,j}, C_{j,i} \subseteq \mathbb{R}$  such that  $u_i(j) \in C_{i,j}$  and  $u_j(i) \in C_{j,i}$  for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$ . Our algorithm uses  $\mathcal{C}$  to get an upper confidence bound for each agent's utility function and then computes a stable matching with transfers as if these upper confidence bounds were the true utilities (see COMPUTEMATCH in Algorithm 2). This can be implemented efficiently if we use, e.g., the Hungarian algorithm [28] to solve (P) and (D).

## 5.2 Main Lemma

The key fact we need to analyze our algorithms is that Subset Instability is upper bounded by the sum of the sizes of the relevant confidence sets, assuming that the confidence sets over the utilities contain the true utilities. (In the following, we again omit the round index  $t$ .)

**LEMMA 5.4.** *Suppose a collection of confidence sets  $\mathcal{C} = \{C_{i,j}, C_{j,i} : (i, j) \in \mathcal{I} \times \mathcal{J}\}$  is such that  $u_i(j) \in C_{i,j}$  and  $u_j(i) \in C_{j,i}$  for all  $(i, j)$ . Then, the instability of  $(X^{\text{UCB}}, \tau^{\text{UCB}}) := \text{COMPUTEMATCH}(\mathcal{C})$  satisfies*

$$I(X^{\text{UCB}}, \tau^{\text{UCB}}; u, \mathcal{A}^t) \leq \sum_{a \in \mathcal{A}} \left( \max(C_{a, \mu_{X^{\text{UCB}}}(a)}) - \min(C_{a, \mu_{X^{\text{UCB}}}(a)}) \right). \quad (6)$$



PROOF. Since  $(X^{\text{UCB}}, \tau^{\text{UCB}})$  is stable with respect to  $u^{\text{UCB}}$ , we have  $I(X^{\text{UCB}}, \tau^{\text{UCB}}; u^{\text{UCB}}, \mathcal{A}^t) = 0$ . Thus, it is equivalent to bound the difference  $I(X^{\text{UCB}}, \tau^{\text{UCB}}; u, \mathcal{A}^t) - I(X^{\text{UCB}}, \tau^{\text{UCB}}; u^{\text{UCB}}, \mathcal{A}^t)$ .

At this stage, it might be tempting to bound this difference using the Lipschitz continuity of Subset Instability (see Proposition 4.4). However, this would only allow us to obtain an upper bound of the form  $\sum_{a \in \mathcal{A}} \max_{a' \in \mathcal{A}} (\max(C_{a,a'}) - \min(C_{a,a'}))$ . The problem with this bound is that it depends on the sizes of the confidence sets for all pairs of agents, including those that are *not* matched in  $X^{\text{UCB}}$ , making it too weak to prove regret bounds for UCB-style algorithms.<sup>10</sup> Thus, we proceed with a more fine-grained analysis.

Define the function

$$f(\mathcal{S}, X, \tau; u) = \left( \max_{X' \in \mathcal{X}_{\mathcal{S}}} \sum_{a \in \mathcal{S}} u_a(\mu_{X'}(a)) \right) - \left( \sum_{a \in \mathcal{S}} u_a(\mu_X(a)) + \tau_a \right).$$

By definition,  $I(X, \tau; u, \mathcal{A}) = \max_{\mathcal{S} \subseteq \mathcal{A}} f(\mathcal{S}, X, \tau; u)$ . It follows that

$$\begin{aligned} I(X^{\text{UCB}}, \tau^{\text{UCB}}; u, \mathcal{A}^t) - I(X^{\text{UCB}}, \tau^{\text{UCB}}; u^{\text{UCB}}, \mathcal{A}^t) \\ \leq \max_{\mathcal{S} \subseteq \mathcal{A}} (f(\mathcal{S}, X^{\text{UCB}}, \tau^{\text{UCB}}; u) - f(\mathcal{S}, X^{\text{UCB}}, \tau^{\text{UCB}}; u^{\text{UCB}})). \end{aligned}$$

To finish, we upper bound  $f(\mathcal{S}, X^{\text{UCB}}, \tau^{\text{UCB}}; u) - f(\mathcal{S}, X^{\text{UCB}}, \tau^{\text{UCB}}; u^{\text{UCB}})$  for each  $\mathcal{S} \subseteq \mathcal{A}$ . We decompose this expression into two terms:

$$\begin{aligned} f(\mathcal{S}, X^{\text{UCB}}, \tau^{\text{UCB}}; u) - f(\mathcal{S}, X^{\text{UCB}}, \tau^{\text{UCB}}; u^{\text{UCB}}) \\ = \underbrace{\left( \max_{X' \in \mathcal{X}_{\mathcal{S}}} \sum_{a \in \mathcal{S}} u_a(\mu_{X'}(a)) - \max_{X' \in \mathcal{X}_{\mathcal{S}}} \sum_{a \in \mathcal{S}} u_a^{\text{UCB}}(\mu_{X'}(a)) \right)}_{(A)} \\ + \underbrace{\left( \sum_{a \in \mathcal{S}} (u_a^{\text{UCB}}(\mu_{X^{\text{UCB}}}(a)) + \tau_a^{\text{UCB}}) - \sum_{a \in \mathcal{S}} (u_a(\mu_{X^{\text{UCB}}}(a)) + \tau_a^{\text{UCB}}) \right)}_{(B)}. \end{aligned}$$

To see that (A) is nonpositive, observe that the maximum weight matching of  $\mathcal{S}$  with respect to  $u$  is no larger than the maximum weight matching of  $\mathcal{S}$  with respect to  $u^{\text{UCB}}$ , since  $u^{\text{UCB}}$  pointwise upper bounds  $u$ . To upper bound (B), observe that the transfers cancel out, so the expression is equivalent to

$$\sum_{a \in \mathcal{S}} (u_a^{\text{UCB}}(\mu_{X^{\text{UCB}}}(a)) - u_a(\mu_{X^{\text{UCB}}}(a))) \leq \sum_{a \in \mathcal{A}} (\max(C_{a, \mu_{X^{\text{UCB}}}(a)}) - \min(C_{a, \mu_{X^{\text{UCB}}}(a)})). \quad \square$$

### 5.3 Instantiations of the Meta-algorithm

As formalized in METAMATCHUCB, the regret bound of Lemma 5.4 suggests a simple approach: At each round, select the matching with transfers returned by COMPUTEMATCH and update confidence sets accordingly. To instantiate METAMATCHUCB, it remains to construct confidence intervals that contain the true utilities with high probability. This last step naturally depends on the assumptions made about the utilities and the noise.

<sup>10</sup>For intuition, consider the classical stochastic multi-armed bandits setting and suppose that we could only guarantee that the loss incurred by an arm is bounded by the maximum of the sizes of the confidence sets over *all* arms. Then, we would only be able to obtain a weak bound on regret, since low-reward arms with large confidence sets may never be pulled.

**ALGORITHM 3:** MATCHUCB: A bandit algorithm for matching with transferable utilities for unstructured preferences.

---

```

1: procedure MATCHUCB( $T$ )
2:   for  $(i, j) \in \mathcal{I} \times \mathcal{J}$  do                                 $\triangleright$  Initialize confidence intervals.
3:      $C_{i,j} \leftarrow [-1, 1]$ 
4:      $C_{j,i} \leftarrow [-1, 1]$ 
5:   for  $1 \leq t \leq T$  do
6:      $(X^t, \tau^t) \leftarrow \text{COMPUTEMATCH}(\mathcal{C})$ 
7:     for  $(i, j) \in X^t$  do                                 $\triangleright$  Set confidence intervals and update means.
8:       Update empirical means  $\hat{u}_i(j)$  and  $\hat{u}_j(i)$  from feedback; increment counter  $n_{ij}$ .
9:        $C_{i,j} \leftarrow \left[ \hat{u}_i(j) - 8\sqrt{\log(|\mathcal{A}|T)/n_{ij}}, \hat{u}_i(j) + 8\sqrt{\log(|\mathcal{A}|T)/n_{ij}} \right] \cap [-1, 1]$ 
10:       $C_{j,i} \leftarrow \left[ \hat{u}_j(i) - 8\sqrt{\log(|\mathcal{A}|T)/n_{ij}}, \hat{u}_j(i) + 8\sqrt{\log(|\mathcal{A}|T)/n_{ij}} \right] \cap [-1, 1]$ 

```

---

**ALGORITHM 4:** MATCHTYPEDUCB: A bandit algorithm for matching with transferable utilities for typed preferences.

---

```

1: procedure MATCHTYPEDUCB( $T$ )
2:   for  $(c, c') \in \mathcal{C} \times \mathcal{C}$  do                                 $\triangleright$  Initialize confidence intervals and empirical means.
3:      $C_{c,c'} \leftarrow [-1, 1]$ 
4:   for  $1 \leq t \leq T$  do
5:      $(X^t, \tau^t) \leftarrow \text{COMPUTEMATCH}(\mathcal{C})$ 
6:     for  $(i, j) \in X^t$  do                                 $\triangleright$  Set confidence intervals and update means.
7:       Update empirical means  $\hat{f}(c_i, c_j)$  and  $\hat{f}(c_j, c_i)$  from feedback; increment  $n_{c_i, c_j}$ .
8:        $C_{c_i, c_j} \leftarrow \left[ \hat{f}(c_i, c_j) - 8\sqrt{\log(|\mathcal{A}|T)/n_{c_i, c_j}}, \hat{f}(c_i, c_j) + 8\sqrt{\log(|\mathcal{A}|T)/n_{c_i, c_j}} \right] \cap [-1, 1]$ 
9:        $C_{c_j, c_i} \leftarrow \left[ \hat{f}(c_j, c_i) - 8\sqrt{\log(|\mathcal{A}|T)/n_{c_i, c_j}}, \hat{f}(c_j, c_i) + 8\sqrt{\log(|\mathcal{A}|T)/n_{c_i, c_j}} \right] \cap [-1, 1]$ 

```

---

*Unstructured preferences.* For this setting, we construct confidence intervals following the classical UCB approach: For each utility value involving the pair  $(i, j) \in \mathcal{I} \times \mathcal{J}$ , we take a confidence interval of length  $O(\sqrt{\log(|\mathcal{A}|T)/n_{ij}})$  centered at the empirical mean, where  $n_{ij}$  is the number of times the pair has been matched thus far. We describe this construction precisely in Algorithm 3 (MATCHUCB).

To analyze MATCHUCB, recall that Lemma 5.4 bounds the regret at each step by the lengths of the confidence intervals of each pair in the selected matching. Bounding the lengths of the confidence intervals parallels the analysis of UCB for classical stochastic multi-armed bandits. We give the full proof of Theorem 5.1 in Appendix C.1.

*Typed Preferences.* For this setting, we construct our confidence intervals as follows: For each pair of types  $c_1$  and  $c_2$ , we take a length  $O(\sqrt{\log(|\mathcal{A}|T)/n_{c_1 c_2}})$  confidence interval centered around the empirical mean, where  $n_{c_1 c_2}$  is the number of times that an agent with type  $c_1$  has been matched with an agent with type  $c_2$ . We describe this construction precisely in Algorithm 4 (MATCHTYPEDUCB). We give the full proof of Theorem 5.2 in Appendix C.2.

*Separable Linear Preferences.* To build the confidence sets, we use a key idea from the design of LinUCB [30, 39]. The idea is to compute a confidence set for each hidden vector  $\phi(a)$  using the least squares estimate and use that to construct confidence sets  $C_{a,a'}$  for the utilities.

**ALGORITHM 5:** MATCHLINUCB: A bandit algorithm for matching with transferable utilities for separable linear preferences.

---

```

1: procedure MATCHLINUCB( $T$ )
2:   for  $(i, j) \in \mathcal{I} \times \mathcal{J}$  do                                ▶ Initialize confidence intervals.
3:      $C_{i,j} \leftarrow [-1, 1]$ 
4:      $C_{j,i} \leftarrow [-1, 1]$ 
5:   for  $1 \leq t \leq T$  do
6:      $(X^t, \tau^t) \leftarrow \text{COMPUTEMATCH}(\mathcal{C})$ 
7:     for  $a \in \mathcal{A}^t$  do                                       ▶ Update confidence intervals.
8:       Increment the counter  $n_a$ .
9:        $\beta \leftarrow O\left(d \log T + \frac{n_a \sqrt{\ln(n_a/(T|A|))}}{T^2}\right)$ .    ▶ Parameter for width of confidence set.
10:      if  $\mu_{X^t}(a) \neq a$  then
11:        Add  $t$  to  $\mathcal{T}_a$  (the set of rounds in which agent  $a$  has been matched).
12:        Set  $\mathcal{R}_{a,t}$  equal to the observed utility for agent  $a$  in round  $t$ .
13:         $\phi^{\text{LS}}(a) \leftarrow \underset{v \in \mathcal{B}^d}{\text{argmin}} \left( \sum_{t' \in \mathcal{T}_a} (\langle v, c_{\mu_{X^{t'}}}(a) \rangle - \mathcal{R}_{a,t'})^2 \right)$     ▶ Least squares
estimate.
14:         $C_{\phi(a)} \leftarrow \left\{ v \mid \sum_{t' \in \mathcal{T}_a} (\langle v - \phi^{\text{LS}}(a), c_{\mu_{X^{t'}}}(a) \rangle)^2 \leq \beta, \|v\|_2 \leq 1 \right\}$     ▶ Conf. ellipsoid.
15:        for  $a' \in \mathcal{A}$  do
16:           $C_{a,a'} \leftarrow \left\{ \langle c_{a'}, v \rangle \mid v \in C_{\phi(a)} \right\} \cap [-1, 1]$     ▶ Update confidence sets of agent  $a$ .
```

---

More formally, let  $\mathcal{T}_a$  be the set of rounds where agent  $a$  is matched on the platform thus far, and for  $t' \in \mathcal{T}_a$ , let  $\mathcal{R}_{a,t'}$  be the observed utility at time  $t'$  for agent  $a$ . The center of the confidence set will be given by the least squares estimate

$$\phi^{\text{LS}}(a) = \arg \min_{v \in \mathcal{B}^d} \left( \sum_{t' \in \mathcal{T}_a} (\langle v, c_{\mu_{X^{t'}}}(a) \rangle - \mathcal{R}_{a,t'})^2 \right).$$

The confidence set for  $\phi(a)$  is given by

$$C_{\phi(a)} := \left\{ v \mid \sum_{t' \in \mathcal{T}_{a,t}} \left\langle v - \phi^{\text{LS}}(a), c_{\mu_{X^{t'}}}(a) \right\rangle^2 \leq \beta \text{ and } \|v\|_2 \leq 1 \right\},$$

where  $\beta = O(D \log T + \frac{n_a \sqrt{\ln(n_a/\delta)}}{T^2})$  and  $n_a$  counts the number of times that  $a$  has appeared in selected matchings. The confidence set for  $u(a, a')$  is given by

$$C_{a,a'} := \left\{ \langle c_{a'}, v \rangle \mid v \in C_{\phi(a)} \right\} \cap [-1, 1].$$

We describe this construction precisely in Algorithm 5 (MATCHLINUCB). We give the full proof of Theorem 5.3 in Appendix C.3.

## 5.4 Matching Lower Bound

For the case of unstructured preferences, we now show that MATCHUCB achieves optimal regret (up to logarithmic factors) by showing a lower bound that (nearly) matches the upper bound in Theorem 5.1.

LEMMA 5.5. *For any algorithm that learns a stable matching with respect to unstructured preferences, there exists an instance on which it has expected regret  $\tilde{\Omega}(|A|^{3/2}\sqrt{T})$  (where regret is given by Subset Instability).*

The idea behind this lemma is to show a lower bound for the easier problem of learning a maximum weight matching using utility difference as regret. By Proposition 4.4, this immediately implies a lower bound for learning a stable matching with regret measured by Subset Instability.

This lower bound illustrates the close connection between our setting and that of learning a maximum weight matching. Indeed, by applying MATCHUCB and simply disregarding the transfers every round, we recover the classical UCB-based algorithm for learning the maximum weight matching [15, 22, 29]. From this perspective, the contribution of MATCHUCB is an approach to set the dual variables while asymptotically maintaining the same regret as the primal-only problem.

## 6 EXTENSIONS

In this section, we discuss several extensions of our results: instance-dependent regret bounds, connections between subset instability and platform revenue, and non-transferable utilities. These extensions illustrate the generality of our framework and also suggest several avenues for future research.

In Section 6.1, we derive instance-dependent regret bounds for Subset Instability, which allow us to improve the  $O(\sqrt{T})$  convergence from Section 5 to  $O(\log T)$  for any fixed instance. Achieving this logarithmic bound involves choosing “robust” dual solutions when setting transfers (rather than choosing an arbitrary optimal primal-dual pair as in COMPUTEMATCH): We want our selected primal-dual pair to lead to stable outcomes even under perturbations of the transfers.

In Section 6.2, we connect the subsidy perspective of Subset Instability to platform revenue. We relate regret to platform revenue and show that, when there are search frictions, the platform can achieve substantial long-run profit despite starting with no knowledge of agent preferences.

In Section 6.3, we adapt our framework to matching with non-transferable utilities (where agents do not transfer money to other agents on the platform). We define an analogue of Subset Instability using the subsidy formulation and give an  $\tilde{O}(\sqrt{T})$  regret algorithm for learning stable matchings.

### 6.1 Instance-dependent Regret Bounds

While our analyses in Section 5.1 focused on bounds that hold uniformly for all problem instances, we now explore *instance-dependent* regret bounds. Instance-dependent bounds capture a different facet of bandit algorithms: How does the number of mistakes made by the algorithm scale on each instance with respect to  $T$ ? Bounds of this nature have been explored in previous works [9, 12, 32, 33, 40] on learning stable matchings in the non-transferable utilities setting, and we show that they can be obtained within our framework as well.

Our instance-dependent regret bound depends on a gap  $\Delta > 0$  determined by the true utility function  $u$ . We focus on the setting where agent utilities are unstructured (i.e.,  $u \in \mathcal{U}_{\text{unstructured}}$ ) and where the same set of agents  $\mathcal{A}$  arrives in each round. As is common in analyses of combinatorial bandit problems (e.g., References [15, 29]), the gap  $\Delta$  in the bound is global to the matching. Letting  $X^{\text{opt}}$  be a maximum weight matching with respect to  $u$ , we define the gap to  $\Delta$  be the difference in utility between the optimal and second-best matchings:<sup>11</sup>

$$\Delta := \inf_{X \neq X^{\text{opt}}} \left\{ \sum_{a \in \mathcal{A}} u_a(\mu_{X^{\text{opt}}}(a)) - \sum_{a \in \mathcal{A}} u_a(\mu_X(a)) \right\}.$$

<sup>11</sup>Our bound is less fine-grained than the gap in Chen et al. [15], and in particular does not allow there to be multiple maximum weight matchings. We defer improving our definition of  $\Delta$  to future work.

We prove the following regret bound:

**THEOREM 6.1 (INSTANCE-DEPENDENT REGRET).** *Suppose that  $\mathcal{A}_t = \mathcal{A}$  for all  $t$ . Let  $u \in \mathcal{U}_{\text{unstructured}}$  be any utility function and put*

$$\Delta := \inf_{X \neq X^*} \left\{ \sum_{a \in \mathcal{A}} u_a(\mu_{X^*}(a)) - \sum_{a \in \mathcal{A}} u_a(\mu_X(a)) \right\}.$$

*Then, MATCHUCB' incurs expected regret  $\mathbb{E}(R_T) = O(|\mathcal{A}|^5 \cdot \log(|\mathcal{A}|T)/\Delta^2)$ .*

*Remark.* MATCHUCB' is MATCHUCB with a slight adjustment to COMPUTEMATCH needed to prove Theorem 6.1. MATCHUCB', like MATCHUCB, does not depend on the gap  $\Delta$  and achieves the instance-independent regret bound in Theorem 5.1.<sup>12</sup> That is, MATCHUCB' achieves both our instance-independent and instance-dependent regret bounds.

Our starting point for proving Theorem 6.1 is to upper bound the number of “mistakes” that a platform makes while exploring and learning, i.e., the number of rounds where the chosen matching is suboptimal. That is, we bound the number of rounds where the chosen market outcome is not stable with respect to the true utilities  $u$ . This is similar in spirit to the analysis of the combinatorial bandits problem of learning a maximum weight matching in Chen et al. [15]. However, a crucial difference is that a mistake can be incurred even when the selected matching is optimal if the selected transfers do not result in a stable market outcome. Ensuring that the selected transfers result in a stable market outcome when the utility estimates are sufficiently accurate is the main technical hurdle in our analysis.

To make this argument work, we need to specify more precisely how the primal-dual solution is chosen in line 5 of COMPUTEMATCH (which we previously did not specify). In particular, poor choices of the primal-dual solution can lead to many rounds where the chosen outcome is unstable, because the transfers violate the stability constraints. To see this, consider a market with a single customer  $C$  and a single provider  $P$  such that  $u_C(P) = 2$  and  $u_P(C) = -1$ , and suppose we have nearly tight upper bounds  $u_C^{\text{UCB}}(P) = 2 + \varepsilon$  and  $u_P^{\text{UCB}}(C) = -1 + \varepsilon$  on the utilities. Then, the market outcome with matching  $\{(C, P)\}$  with  $\tau_C = -2 - \varepsilon$  and  $\tau_P = -\tau_C$  could be selected by COMPUTEMATCH, since it corresponds to an optimal primal-dual pair for  $u^{\text{UCB}}$ . However, it is not stable with respect to the true utilities  $u$  (as individual rationality is violated for  $C$ ), regardless of how small  $\varepsilon$  is. Thus, without assuming more about how the optimal primal-dual pair is chosen in COMPUTEMATCH, we cannot hope to bound the number of unstable market outcomes selected.

We show that, by carefully selecting an optimal primal-dual pair each round, we can bound the number of mistakes. In particular, we design an algorithm COMPUTEMATCH' to find primal-dual pairs that satisfy the following property: If the confidence sets are small enough, then the selected matching will be stable with respect to the true utilities.

**LEMMA 6.2.** *Suppose COMPUTEMATCH' is run on a collection  $\mathcal{C}$  of confidence sets  $C_{i,j}$  and  $C_{j,i}$  over the agent utilities that satisfy*

$$\max(C_{i,j}) - \min(C_{i,j}) \leq 0.05 \frac{\Delta}{|\mathcal{A}|} \quad \text{and} \quad \max(C_{j,i}) - \min(C_{j,i}) \leq 0.05 \frac{\Delta}{|\mathcal{A}|}$$

*for all  $(i,j)$  in the matching returned by COMPUTEMATCH'. Suppose also that the confidence sets  $\mathcal{C}$  contain the true utilities for all pairs of agents. Then, the market outcome returned by COMPUTEMATCH' is stable with respect to the true utilities  $u$ .*

<sup>12</sup>The instance-independent regret bound can be shown using the same argument as the proof for Theorem 5.1.

*Remark.* Lemma 6.2 does *not* hold for COMPUTEMATCH; its proof relies on the particular specification of the optimal primal-dual pair in COMPUTEMATCH'.

Using Lemma 6.2, we intuitively can bound the number of mistakes made by the algorithm by the number of samples needed to sufficiently reduce the size of the confidence sets. In Appendix D, we describe how we choose optimal primal-dual pairs in COMPUTEMATCH', prove Lemma 6.2, and provide a full proof of Theorem 6.1.

Theorem 6.1 opens the door to further exploring algorithmic properties of learning stable matchings. First, this result establishes fine-grained regret bounds, demonstrating the typical  $O(\log T)$  regret bounds from the combinatorial bandits literature [15] are achievable in our setting as well. Second, Theorem 6.1 provides insight into the number of mistakes made by the platform. In particular, we show within the proof of Theorem 6.1 that the platform fails to choose a matching that is stable with respect to  $u$  in at most  $O(|\mathcal{A}|^4 \cdot \log(|\mathcal{A}|T)/\Delta^2)$  rounds.<sup>13</sup> This means that the platform selects a stable matching in at least  $T - O(|\mathcal{A}|^5 \cdot \log(|\mathcal{A}|T)/\Delta^2) = T - O(\log T)$  of the rounds.

As we described, our bounds in Theorem 6.1 rely on choosing an appropriate primal-dual solution. An interesting direction for future work would be to provide further insight into how different methods for finding optimal primal-dual pairs affect both regret bounds and the trajectory of the selected market outcomes over time.

## 6.2 Search Frictions and Platform Revenue

Next, we further ground Subset Instability by explicitly connecting it to the platform's revenue under a stylized economic model of search frictions. A major motivation for this is that it helps explain when an online platform can earn a profit in competitive settings, even when they start out with no information about agent preferences.

More specifically, we incorporate search frictions where an agent must lose utility  $\varepsilon$  to find an alternative to the given match (e.g., from the time spent finding an alternate partner or from a cancellation fee). These search frictions weaken the requirements for stability: The platform now only needs matchings to be  $\varepsilon$ -stable:

$$u_i(j) + u_j(i) - 2\varepsilon \leq u_i(\mu_X(i)) + \tau_i + u_j(\mu_X(j)) + \tau_j$$

for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$  and  $u_a(\mu_X(a)) + \tau_a \geq -\varepsilon$  for all  $a \in \mathcal{A}$ .<sup>14</sup>

To model revenue, we take the subsidy perspective on Subset Instability. Specifically, recall that Subset Instability is equal to the minimum subsidy needed to maintain stability (see Proposition 4.2). With search frictions, that subsidy can potentially be *negative*, thus allowing the platform to generate revenue. We are interested in analyzing the maximum revenue (minimum subsidy) the platform can generate while ensuring stability with high probability over all rounds. For realism, we also want this subsidy to be computed online using only information that the platform has access to, but it turns out we can do this with minimal modifications to our algorithm.

More formally, in this modified model, the platform must select an  $\varepsilon$ -stable matching in each round with high probability by choosing appropriate subsidies. That is, in round  $t$ , the platform

<sup>13</sup>The number of mistakes necessarily depends on the gap  $\Delta$ , because there exist utility functions  $u$  and  $\tilde{u}$  where  $\|u - \tilde{u}\|_\infty$  is arbitrarily small, but where the stable market outcomes with respect to  $u$  and  $\tilde{u}$  differ. To see this, consider a market where  $\mathcal{I} = \{C\}$  and  $\mathcal{J} = \{P\}$ . Suppose that  $u_C(P) = \tilde{u}_C(P) = 1$ , while  $u_P(C) = -1 + \varepsilon$  and  $\tilde{u}_P(C) = -1 - \varepsilon$ . Then, the maximum weight matchings under these utility functions differ:  $\{(C, P)\}$  is the only maximum weight matching in the former, whereas  $\emptyset$  is the only maximum weight matching in the latter.

<sup>14</sup>This definition corresponds to  $(X, \tau)$  belonging to the weak  $\varepsilon$ -core of Shapley and Shubik [41]. We note that this definition also relaxes individual rationality. This formulation gives us the cleanest algorithmic results; while it can be extended to an analogue that does not relax individual rationality, it would involve bounds that (necessarily) depend on the specifics of agents' utilities.



selects a matching with transfers  $(X^t, \tau^t)$  with the modification that the transfers need not be zero-sum. The transfers thus incorporate the amount that platform is subsidizing or charging agents for participation on the platform. The net profit of the platform is then  $-\sum_{t=1}^T \sum_{a \in \mathcal{A}} \tau_a^t$ . We impose the stability requirement that

$$\mathbb{P}[(X^t, \tau^t) \text{ is } \varepsilon\text{-stable for all } 1 \leq t \leq T] \geq 0.99.$$

Given this setup, we show the following:

**THEOREM 6.3.** *For preference class  $\mathcal{U}_{\text{unstructured}}$  (see Section 3), there exists an algorithm giving the platform*

$$\varepsilon T \sum_{t=1}^T |\mathcal{A}_t| - O\left(|\mathcal{A}| \sqrt{nT} \sqrt{\log(|\mathcal{A}|T)}\right)$$

*revenue in the presence of search frictions while maintaining stability with high probability.*

*Remark.* In particular, if  $\mathcal{A}_t = \mathcal{A}$  in every round, then the platform will start making a profit within  $O(|\mathcal{A}|/\varepsilon^2 \cdot \log(|\mathcal{A}|/\varepsilon^2))$  rounds.

We defer the proof of Theorem 6.3 to Appendix E.

Qualitatively, Theorem 6.3 captures that if the platform “pays to learn” in initial rounds, then the information that it obtains will help it achieve a profit in the long run. We note that both the revenue objective and the model for search frictions that we consider in these preliminary results are stylized. An interesting direction for future work would be to integrate more realistic platform objectives and models for search frictions into the framework.

### 6.3 Matching with Non-transferable Utilities

While we have focused on matching with transferable utilities, utilities are not always transferable in practice, as in the cases of dating markets and college admissions (i.e., most people are not willing to date an undesirable partner in exchange for money, and a typical college admission slot is not sold for money). We can extend our findings to this setting following the model of matching with **non-transferable utilities (NTU)** [23], which has also been studied in previous work [12, 17, 32, 40]. The definition of Subset Instability extends naturally and has advantages over the “utility difference” metric that is commonly used in prior work. Our algorithmic meta-approach also sheds new light on the convergence properties of the centralized UCB algorithm of Liu et al. [32].

The starting point of our instability measure is slightly different than in Section 4. Since stable matchings in the NTU model need not maximize total utility, we cannot define instability based on a maximum over all subsets of agents of the utility difference for that subset. However, the subsidy formulation of Subset Instability (see (4)) translates well to this setting. Our instability measure will correspond to the minimum amount the platform could subsidize agents so individual rationality holds and no blocking pairs remain. For matching with NTU, we formalize this notion as follows:

**Definition 6.4 (NTU Subset Instability).** For utilities  $u$  and agents  $\mathcal{A}$ , the *NTU Subset Instability*  $I(X; u, \mathcal{A})$  of a matching  $X$  is

$$\begin{aligned} \min_{s \in \mathbb{R}^{|\mathcal{A}|}} \sum_{a \in \mathcal{A}} s_a & \tag{†} \\ \text{s.t. } \min(u_i(j) - u_i(\mu_X(i)) - s_i, u_j(i) - u_j(\mu_X(j)) - s_j) & \leq 0 \quad \forall (i, j) \in \mathcal{I} \times \mathcal{J} \\ u_a(\mu_X(a)) + s_a & \geq 0 \quad \forall a \in \mathcal{A} \\ s_a & \geq 0 \quad \forall a \in \mathcal{A}. \end{aligned}$$

NTU Subsidy Instability inherits some of the same appealing properties as Subsidy Instability.

PROPOSITION 6.5 (INFORMAL). *NTU Subset Instability satisfies the following properties:*

- (1) *NTU Subset Instability is always nonnegative and is zero if and only if  $(X, \tau)$  is stable.*
- (2) *NTU Subset Instability is Lipschitz continuous with respect to agent utilities. That is, for any matching  $X$  and any pair of utility functions  $u$  and  $\tilde{u}$ , it holds that:*

$$|I(X; u, \mathcal{A}) - I(X; \tilde{u}, \mathcal{A})| \leq 2 \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_{\infty}.$$

The proofs of this and subsequent results are deferred to Appendix F. Together, the preceding properties mean that NTU Subsidy Instability is useful as a regret measure for learning stable matchings.

As in the transferable utilities setting, Property 2 implies the existence of an explore-then-commit algorithm with  $\tilde{O}(|\mathcal{A}|^{4/3} T^{2/3})$  regret. We show that this can be improved to a  $\sqrt{T}$  dependence by adapting our approach from Section 5:

THEOREM 6.6. *For matchings with non-transferable utilities, there exists an algorithm that for any utility function  $u$  incurs regret  $R_T = O(|\mathcal{A}|^{3/2} \sqrt{T} \sqrt{\log(|\mathcal{A}|T)})$ .*

While Theorem 6.6 illustrates that our approach easily generalizes to the NTU setting, we highlight two crucial differences between these settings. First, learning a stable matching is incomparable to learning a maximum weight matching, because stable matchings do not maximize the sum of agents' utilities in the NTU setting. Next, the instability measure is *not* equivalent to the cumulative unhappiness of agents, unlike in the setting with transferable utilities. Intuitively, these definitions cease to be equivalent, because non-transferable utilities render the problem more “discontinuous” and thus obstruct the duality results we applied earlier.

These results provide a preliminary application of our framework to the setting of matching with non-transferable utilities; an interesting direction for future inquiry would be to more thoroughly investigate notions of approximate stability and regret in this setting.

**6.3.1 Comparison to the Utility Difference Measure.** It turns out that the algorithm underlying Theorem 6.6 is equivalent to the centralized UCB algorithm from previous work [12, 32], albeit derived from a different angle. However, an important difference is that Theorem 6.6 guarantees low regret relative to the incentive-aware NTU Subset Instability, as opposed to the incentive-unaware “utility difference” measure in prior work. In this section, we outline several properties that make our instability measure more suitable especially in the NTU setting. In particular, we show for utility difference that:

- (a) There is no canonical formalization of utility difference when multiple stable matchings exist.
- (b) The utility difference of a matching can be positive even if the matching is stable and negative even if the matching is unstable.
- (c) Even when restricting to markets with unique stable matchings, the utility difference of a matching can be discontinuous in the true agent utilities. As a result, it does not allow for instance-independent regret bounds that are sublinear in  $T$ .

For (a), the utility difference requires specifying a stable matching to serve as a benchmark against which to measure relative utility. However, when multiple stable matchings exist, some ambiguity arises as to which one should be chosen as the benchmark. Because of this, previous works [12, 17, 32, 40] study two different benchmarks. In particular, they assume providers' preferences are known and benchmark with respect to the customer-optimal and customer-pessimal

stable matchings. For (b), notice that the utility difference for the maximum weight matching is negative, even though it is typically not stable in the NTU setting. Moreover, because of the ambiguity in the benchmark from (a), the utility difference may not be zero even when the matching is stable. For (c), to see that utility difference is not continuous as a function of the underlying agent utilities, consider the following example:

*Example 6.7.* Consider a market where there is a single customer  $i$  and two providers  $j_1$  and  $j_2$ . Suppose their utility functions are given by  $u_i(j_1) = \varepsilon$ ,  $u_i(j_2) = 2\varepsilon$ ,  $u_{j_1}(i) = 1$ , and  $u_{j_2}(i) = 0.5$ . Then, the unique stable matching  $\{(i, j_2)\}$  has total utility  $0.5 + 2\varepsilon$ . Now, consider the perturbed utility function  $\tilde{u}$  such that  $\tilde{u}_i(j_1) = 2\varepsilon$ ,  $\tilde{u}_i(j_2) = \varepsilon$ ,  $\tilde{u}_{j_1}(i) = 1$ , and  $\tilde{u}_{j_2}(i) = 0.5$ . For this perturbed utility function, the unique stable matching is  $\{(i, j_1)\}$ , which has total utility  $1 + 2\varepsilon$ . The utility difference (either optimal or pessimal) for matching  $\{(i, j_2)\}$  is 0 for  $u$  and  $0.5 + \varepsilon$  for  $\tilde{u}$ . Since this holds for any  $\varepsilon > 0$ , taking  $\varepsilon \rightarrow 0$  shows that utility difference is not continuous in the utility function.

That utility difference is discontinuous in agent utilities rules out the existence of bandit algorithms that achieve sublinear instance-independent regret when using utility difference as the regret measure. In particular, the analyses in previous work [12, 32, 33, 40] focus entirely on *instance-dependent* regret bounds. They show that centralized UCB achieves logarithmic instance-dependent regret with respect to the utility difference relative to the customer-pessimal stable matching (but does not achieve sublinear regret with respect to the customer-optimal stable matching). Our insight here is that a new measure of instability can present a more appealing evaluation metric and paint a clearer picture of an algorithm's convergence to the *set* of stable matchings as a whole.

## 7 IN WHAT SETTINGS ARE EQUILIBRIA LEARNABLE?

A core insight of our work is that, in a stochastic environment, “optimism in the face of uncertainty” can be effectively leveraged for the problem of learning stable matchings. This motivates us to ask: In what other settings, and with what other algorithmic methods, can equilibria be learned?

One interesting open direction is to understand when equilibria can be learned in *adversarial* environments where the utility functions can change between rounds. From an economic perspective, adversarial environments could capture evolving market conditions. In the adversarial bandit setting, most work relies on gradient-based algorithms instead of UCB-based algorithms to attain optimal regret bounds (see, e.g., References [1, 6]). Can these gradient-based algorithms similarly be adapted to Subset Instability?

Another interesting open direction is to consider more general market settings, even in stochastic environments. For example, within the context of matching markets, each agent might match to more than one agent on the other side of the market; and outside of matching markets, a buyer might purchase multiple units of multiple goods. In markets with transferable utilities, incentive-aligned outcomes can be captured by *Walrasian equilibria* (see, e.g., Bichler et al. [10]). Can Subset Instability and our UCB-based algorithms be adapted to learning Walrasian equilibria in general?

Addressing these questions would provide a richer understanding of when and how large-scale, data-driven marketplaces can efficiently learn market equilibria.

## APPENDICES

### A CLASSICAL RESULTS FOR MATCHING WITH TRANSFERABLE UTILITIES

To be self-contained, we briefly state and prove the key results from Shapley and Shubik [42] we need.

First, we explicitly relate the primal-dual formulation in Section 5.1 to stable matchings.

**THEOREM A.1** ([42]). *If  $(X, \tau)$  is stable, then  $(Z, p)$  is an optimal primal-dual pair to (P) and (D), where  $p_a = \tau_a + u_a(X(a))$  and  $Z$  is the indicator matrix in  $\mathbb{R}^{I \times J}$  corresponding to  $X$ .*

*Moreover, if  $(Z, p)$  is an optimal primal-dual pair to (P) and (D) such that  $Z$  lies at an extreme point of the feasible set, then  $(X, \tau)$  is stable where  $\tau_a = p_a - u_a(X(a))$  and  $X$  is the matching corresponding to the nonzero entries of  $Z$ .*

**PROOF.** Both statements follow from the complementary slackness conditions and the definition of stability in Definition 2.1. The complementary slackness conditions are:

- If  $Z_{i,j} > 0$ , then  $p_i + p_j = u_i(j) + u_j(i)$ .
- If  $p_i > 0$ , then  $\sum_j Z_{i,j} = 1$ .
- If  $p_j > 0$ , then  $\sum_i Z_{i,j} = 1$ .

Suppose that  $(X, \tau)$  is stable. Let us first show that  $(Z, p)$  is feasible. We see that  $Z$  is primal feasible by definition. For dual feasibility, since there are no blocking pairs, we know that

$$(u_i(\mu_X(i)) + \tau_i) + (u_j(\mu_X(j)) + \tau_j) \geq u_i(j) + u_j(i),$$

which implies

$$p_i + p_j \geq u_i(j) + u_j(i).$$

The individual rationality condition  $u_a(\mu_X(a)) + \tau_a \geq 0$  tells us  $p_a \geq 0$ . Hence,  $p$  is dual feasible. Next, we show that  $(Z, p)$  is an optimal primal-dual pair by checking the Karush–Kuhn–Tucker conditions. We have already shown primal and dual feasibility, so it suffices to show complementary slackness. The first condition follows from zero-sum transfers. To see the second and third conditions, we show the contrapositive: If  $i \in I$  is such that  $\sum_j Z_{i,j} < 1$ , then  $\sum_j Z_{i,j} = 0$  by our assumption on  $Z$ . Hence,  $i$  is unmatched (i.e.,  $u_i(\mu_X(i)) = 0$  and  $\tau_i = 0$ ), which implies  $p_i = 0$ . The analogous argument applies for  $j \in J$ .

We now prove the second part of the theorem. Suppose  $(Z, p)$  is an optimal solution to (P) and (D) such that  $Z$  is at a vertex. By the Birkhoff–von Neumann theorem, since  $Z$  is a vertex, it corresponds to a matching. We wish to show that  $(X, \tau)$  has no blocking pairs, is individually rational, and has zero-sum transfers. Dual feasibility tells us that:

$$p_i + p_j \geq u_i(j) + u_j(i),$$

which means that:

$$(u_i(\mu_X(i)) + \tau_i) + (u_j(\mu_X(j)) + \tau_j) \geq u_i(j) + u_j(i),$$

so there are no blocking pairs. Dual feasibility also tells us that  $p_a \geq 0$ , which means that  $u_a(\mu_X(a)) + \tau_a \geq 0$ , so individual rationality is satisfied. To show that there are zero-sum transfers, we use complementary slackness. The first complementary slackness condition tells us that if  $Z_{i,j} > 0$ , then  $p_i + p_j = u_i(j) + u_j(i)$ . Using the fact that  $Z$  corresponds to a matching, this in particular means that if  $(i, j) \in X$ , then we know  $\tau_i + \tau_j = 0$ . To show that agents who are unmatched receive 0 transfers, let us use the second and third complementary slackness conditions. The contrapositive tells us that if  $a$  is unmatched, then  $p_a = 0$ , which implies  $\tau_a = 0$ .  $\square$

Since (P) is exactly the maximum weight matching linear program, Theorem A.1 immediately tells us that if  $(X, p)$  is stable, then  $X$  is a maximum weight matching. This means that stable matchings with transferable utilities maximize social welfare.

## B PROOFS FOR SECTION 4

This section contains further exposition (including proofs) for Section 4.

### B.1 Limitations of Utility Difference as an Instability Measure

To illustrate why utility difference fails to be a good measure of instability, we describe a matching with transfers that (i) is far from stable, and (ii) has zero utility difference (but large Subset Instability).

*Example B.1.* Consider the following market with two agents:  $\mathcal{I} = \{i\}$  and  $\mathcal{J} = \{j\}$ . Suppose that  $u_i(j) = 2$  and  $u_j(i) = -1$ . Consider the matching  $X = \{(i, j)\}$  with transfers  $\tau_i = -\xi$  and  $\tau_j = \xi$  for some  $\xi > 0$ . We will show that this matching with transfers will have the properties stated above when  $\xi$  is large.

This matching with transfers has a utility difference equal to zero (for any  $\xi$ ), since it maximizes the sum of utilities. Indeed, it is stable for any  $\xi \in [1, 2]$ . However, when  $\xi > 2$ , this matching with transfers is no longer stable, since the individual rationality condition  $u_i(j) + \tau_i \geq 0$  fails. (Intuitively, the larger  $\xi$  is, the further we are from stability.) But its utility difference remains at zero.

However, the Subset Instability of this matching with transfers is  $\xi - 2 > 0$  when  $\xi > 2$ . In particular, Subset Instability increases with  $\xi$  in this regime, which is consistent with the intuition that outcomes with larger  $\xi$  should be more unstable.

### B.2 Proof of Proposition 4.2

PROPOSITION 4.2. *Minimum stabilizing subsidy equals Subset Instability for any market outcome.*

PROOF OF PROPOSITION 4.2. We can take the dual of the linear program (5) to obtain:

$$\begin{aligned} \max_{\substack{S \in \mathbb{R}^{|\mathcal{I}| \times |\mathcal{J}|} \\ Z \in \mathbb{R}^{|\mathcal{A}|}}} \quad & \sum_{(i,j) \in \mathcal{I} \times \mathcal{J}} S_{i,j} \left( (u_i(j) - u_i(\mu_X(i)) - \tau_i) + (u_j(i) - u_j(\mu_X(j)) - \tau_j) \right) \\ & - \sum_{a \in \mathcal{A}} Z_a (u_a(\mu_X(a)) + \tau_a) \quad (\ddagger) \\ \text{s.t.} \quad & Z_i + \sum_{j \in \mathcal{J}} S_{i,j} \leq 1 \quad \forall i \in \mathcal{I}; \quad Z_j + \sum_{i \in \mathcal{I}} S_{i,j} \leq 1 \quad \forall j \in \mathcal{J}; \\ & S_{i,j} \geq 0 \quad \forall (i,j) \in \mathcal{I} \times \mathcal{J}; \quad Z_a \geq 0 \quad \forall a \in \mathcal{A}. \end{aligned}$$

By strong duality, the optimal values of (5) and  $(\ddagger)$  are equal. Thus, it suffices to show that Subset Instability is equal to  $(\ddagger)$ . By Proposition 4.3, we know that Subset Instability is equal to the maximum unhappiness of any coalition. Thus, it suffices to show that  $(\ddagger)$  is equal to the maximum unhappiness of any coalition.

To interpret  $(\ddagger)$ , observe that there exist optimal  $S^*$  and  $Z^*$  all of whose entries lie in  $\{0, 1\}$ , because this linear program can be embedded into a maximum weight matching linear program. Take such a choice of optimal  $S^*$  and  $Z^*$ . Then,  $S^*$  is an indicator vector corresponding to a (partial) matching on a subset of the agents such that all pairs in this matching are blocking with respect to  $(X, \tau)$ . Similarly,  $Z^*$  is an indicator vector of agents who would rather be unmatched than match according to  $(X, \tau)$ .

We first prove the claim that  $I(X, \tau; u, \mathcal{A})$  is at least  $(\ddagger)$ . Based on the above discussion, the optimal objective of  $(\ddagger)$  is obtained through  $S^*$  and  $Z^*$  that represent a matching and a subset of agents, respectively. Let  $S$  be the union of agents participating in  $S^*$  and  $Z^*$ . We see that the objective of  $(\ddagger)$  is equal to the utility difference at  $S$ , i.e.:

$$\left( \max_{X' \in \mathcal{X}_S} \sum_{a \in S} u_a(\mu_{X'}(a)) \right) - \left( \sum_{a \in S} u_a(\mu_X(a)) + \tau_a \right).$$

This is no larger than Subset Instability by definition.

We next prove the claim that  $I(X, \tau; u, \mathcal{A})$  is at most  $(\ddagger)$ . Let us consider  $S^*$  that maximizes:

$$\max_{S \subseteq \mathcal{A}} \left( \max_{X' \in \mathcal{X}_S} \sum_{a \in S} u_a(\mu_{X'}(a)) \right) - \left( \sum_{a \in S} u_a(\mu_X(a)) + \tau_a \right).$$

Let us take the maximum weight matching of  $S^*$ . Let  $S$  be given by the matched agents in this matching and let  $Z$  be given by the unmatched agents in this matching (using the interpretation of  $(\ddagger)$  described above). We see that the objective at  $(\ddagger)$  for  $(S, Z)$  is equal to Subset Instability, which proves the desired statement.  $\square$

### B.3 Proof of Proposition 4.3

We first formally define the *unhappiness of a coalition*, as follows. In particular, the unhappiness with respect to  $(X, \tau)$  of a coalition  $\mathcal{S} \subseteq \mathcal{A}$  is defined to be:

$$\begin{aligned} & \sup_{\substack{X' \in \mathcal{X}_S \\ \tau' \in \mathbb{R}^{|\mathcal{S}|}}} \sum_{a \in S} (u_a(\mu_{X'}(a)) + \tau'_a) - \sum_{a \in S} (u_a(\mu_X(a)) + \tau_a) \\ & \text{s.t. } u_a(\mu_{X'}(a)) + \tau'_a \geq u_a(\mu_X(a)) + \tau_a \quad \forall a \in S \\ & \quad \tau'_a + \tau'_{\mu_{X'}(a)} = 0 \quad \forall a \in S, \end{aligned} \quad (7)$$

with unhappiness being 0 if there are no feasible  $X'$  and  $\tau'$ . In the optimization program,  $(X', \tau')$  represents a matching with transfers over  $\mathcal{S}$ , with the constraint  $\tau'_a + \tau'_{\mu_{X'}(a)} = 0$  ensuring that it is zero-sum. The objective measures the difference between  $(X, \tau)$  and  $(X', \tau')$  of the total utility of the agents in  $\mathcal{S}$ . The constraint  $u_a(\mu_{X'}(a)) + \tau'_a \geq u_a(\mu_X(a)) + \tau_a$  encodes the requirement that all agents be at least as well off under  $(X', \tau')$  as they were under  $(X, \tau)$ . This optimization program therefore captures the objective of  $\mathcal{S}$  to maximize their total payoff while ensuring that no member of the coalition is worse off than they were according to  $(X, \tau)$ .

Recall that, in terms of unhappiness, Proposition 4.3 is as follows:

**PROPOSITION 4.3.** *The maximum unhappiness of any coalition  $\mathcal{S} \subseteq \mathcal{A}$  with respect to  $(X, \tau)$  equals the Subset Instability  $I(X, \tau; u, \mathcal{A})$ .*

**PROOF OF PROPOSITION 4.3.** By Proposition 4.2, we know that Subset Instability is equal to (5). Moreover, by strong duality, we know that Subset Instability is equal to  $(\ddagger)$  (the dual linear program of (5)). Thus, it suffices to prove that the maximum unhappiness of any coalition is equal to  $(\ddagger)$ .

We first prove the claim that  $(\ddagger)$  is at most the maximum unhappiness of any coalition with respect to  $(X, \tau)$ . To do this, it suffices to construct a coalition  $\mathcal{S} \subseteq \mathcal{A}$  such that  $(\ddagger)$  is at most the unhappiness of  $\mathcal{S}$ . We construct  $\mathcal{S}$  as follows: Recall that there exist optimal solutions  $S^*$  and  $Z^*$  to  $(\ddagger)$  such that  $S^*$  corresponds to a (partial) matching on  $\mathcal{I} \times \mathcal{J}$  and  $Z^*$  corresponds to a subset of  $\mathcal{A}$ . We may take  $\mathcal{S}$  to be the union of the agents involved in  $S^*$  and in  $Z^*$ . Now, we upper bound the unhappiness of  $\mathcal{S}$  by constructing  $X'$  and  $\tau'$  that are feasible for (7). We can take  $X'$  to be the matching that corresponds to the indicator vector  $S^*$ . Because  $(S^*, Z^*)$  is optimal for  $(\ddagger)$ ,

$$u_i(j) + u_j(i) \geq (u_i(\mu_X(i)) + \tau_i) + (u_j(\mu_X(j)) + \tau_j)$$

for all  $(i, j) \in X'$ . Thus, we can find a vector  $\tau'$  of transfers that is feasible for (7). Then, since  $\sum_{a \in S} \tau'_a = 0$ , the objective of (7) at  $(X', \tau')$  is

$$\sum_{a \in S} (u_a(\mu_{X'}(a)) - u_a(\mu_X(a)) - \tau_a).$$

This equals to the objective of  $(\ddagger)$  at  $(S^*, Z^*)$ , which equals  $(\ddagger)$ , as desired.



We now show the inequality in the other direction, that  $(\ddagger)$  is at least the maximum unhappiness of any coalition with respect to  $(X, \tau)$ . It suffices to construct a feasible solution  $(S, Z)$  to  $(\ddagger)$  that achieves at least the maximum unhappiness of any coalition. Let  $\mathcal{S}$  be a coalition with maximum unhappiness, and let  $(X', \tau')$  be an optimal solution for (7). Moreover, let  $S$  be the indicator vector corresponding to agents who are matched in  $X'$  and  $Z$  be the indicator vector corresponding to agents in  $\mathcal{S}$  who are unmatched. The objective of (7) at  $(X', \tau')$  is

$$\sum_{a \in \mathcal{S}} (u_a(\mu_{X'}(a)) - u_a(\mu_X(a)) - \tau_a),$$

which equals the objective of  $(\ddagger)$  at the  $(S, Z)$  that we constructed.  $\square$

#### B.4 Proof of Proposition 4.4

PROPOSITION 4.4. *Subset Instability satisfies the following properties:*

- (1) *Subset Instability is always nonnegative and is zero if and only if  $(X, \tau)$  is stable.*
- (2) *Subset Instability is Lipschitz continuous with respect to agent utilities. That is, for any possible market outcome  $(X, \tau)$  and any pair of utility functions  $u$  and  $\tilde{u}$  it holds that:*

$$|I(X, \tau; u, \mathcal{A}) - I(X, \tau; \tilde{u}, \mathcal{A})| \leq 2 \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_{\infty}.$$

- (3) *Subset Instability is always at least the utility difference.*

PROOF OF PROPOSITION 4.4. We first prove the third part of the Proposition statement, then the first part of the Proposition statement, and finally the second part.

*Proof of part (c).* Because  $\sum_{a \in \mathcal{A}} \tau_a = 0$ , Subset Instability satisfies the following:

$$\begin{aligned} I(X, \tau; u, \mathcal{A}) &\geq \left( \max_{X' \in \mathcal{X}_{\mathcal{A}}} \sum_{a \in \mathcal{A}} u_a(\mu_{X'}(a)) \right) - \left( \sum_{a \in \mathcal{A}} u_a(\mu_X(a)) + \tau_a \right) \\ &= \left( \max_{X' \in \mathcal{X}_{\mathcal{A}}} \sum_{a \in \mathcal{A}} u_a(\mu_{X'}(a)) \right) - \left( \sum_{a \in \mathcal{A}} u_a(\mu_X(a)) \right). \end{aligned}$$

The second line is exactly the utility difference.

*Proof of part (a).* From above, we have that Subset Instability is lower bounded by the utility difference, which is always nonnegative. Hence, Subset Instability is also always nonnegative.

To see that Subset Instability is 0 if and only if  $(X, \tau)$  is stable, first suppose  $(X, \tau)$  is unstable. Then, there exists a blocking pair  $(i, j)$ , in which case

$$I(X, \tau; u, \mathcal{A}) \geq u_i(j) + u_j(i) - (u_i(\mu_X(i)) + u_j(\mu_X(j)) + \tau_i + \tau_j) > 0$$

by the definition of blocking. Now, suppose  $I(X, \tau; u, \mathcal{A}) > 0$ . Then, there exists a subset  $S \subseteq \mathcal{A}$  such that

$$\left( \max_{X' \in \mathcal{X}_S} \sum_{a \in S} u_a(\mu_{X'}(a)) \right) - \left( \sum_{a \in S} u_a(\mu_X(a)) + \tau_a \right) > 0.$$

Let  $X'$  be a maximum weight matching on  $S$ . We can rewrite the above as

$$\sum_{(i,j) \in X'} (u_i(j) + u_j(i) - (u_i(\mu_X(i)) + u_j(\mu_X(j)) + \tau_i + \tau_j)) > 0.$$

Some term in the sum on the left-hand side must be positive, so there exists a blocking pair  $(i, j) \in X'$ . In particular,  $(X, \tau)$  is not stable.

*Proof of part (b).* We prove that

$$|I(X, \tau; u, \mathcal{A}) - I(X, \tau; \tilde{u}, \mathcal{A})| \leq 2 \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_\infty.$$

The supremum of  $L$ -Lipschitz functions is  $L$ -Lipschitz, so it suffices to show that

$$\left( \max_{X' \in \mathcal{X}_S} \sum_{a \in S} u_a(\mu_{X'}(a)) \right) - \sum_{a \in S} (u_a(\mu_X(a)) + \tau_a)$$

satisfies the desired Lipschitz condition for any  $S \subseteq \mathcal{A}$ . In particular, it suffices to show that

$$\left| \sum_{a \in S} (u_a(\mu_X(a)) + \tau_a) - \sum_{a \in S} (\tilde{u}_a(\mu_X(a)) + \tau_a) \right| \leq \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_\infty \quad (8)$$

and

$$\left| \left( \max_{X' \in \mathcal{X}_S} \sum_{a \in S} u_a(\mu_{X'}(a)) \right) - \left( \max_{X' \in \mathcal{X}_S} \sum_{a \in S} \tilde{u}_a(\mu_{X'}(a)) \right) \right| \leq \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_\infty. \quad (9)$$

For (8), we have

$$\left| \sum_{a \in S} (u_a(\mu_X(a)) + \tau_a) - \sum_{a \in S} (\tilde{u}_a(\mu_X(a)) + \tau_a) \right| = \left| \sum_{a \in S} (u_a(\mu_X(a)) - \tilde{u}_a(\mu_X(a))) \right| \leq \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_\infty.$$

For (9), this boils down to showing that the total utility of the maximum weight matching is Lipschitz. Using again the fact that the supremum of Lipschitz functions is Lipschitz, this follows from the total utility of any fixed matching being Lipschitz.  $\square$

## C PROOFS FOR SECTION 5

### C.1 Proof of Theorem 5.1

**THEOREM 5.1.** *For preference class  $\mathcal{U}_{\text{unstructured}}$  (see Section 3), MATCHUCB (defined in Section 5.3) incurs expected regret  $\mathbb{E}(R_T) = O(|\mathcal{A}| \sqrt{nT \log(|\mathcal{A}|T)})$ , where  $n = \max_t |\mathcal{A}_t|$ .*

**PROOF OF THEOREM 5.1.** The starting point for our proof of Theorem 5.1 is the typical approach in multi-armed bandits and combinatorial bandits [15, 22, 30] of bounding regret in terms of the sizes of the confidence interval of the chosen arms. However, rather than using the sizes of confidence intervals to bound the utility difference (as in the incentive-free maximum weight matching setting), we bound Subset Instability through Lemma 5.4. From here on, our approach composes cleanly with existing bandits analyses; in particular, we can follow the typical combinatorial bandits approach [15, 22] to get the desired upper bound.

For completeness, we present the full proof. We divide into two cases, based on the event  $E$  that all of the confidence sets contain their respective true utilities at every timestep  $t \leq T$ . That is,  $u_i(j) \in C_{i,j}$  and  $u_j(i) \in C_{j,i}$  for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$  at all  $t$ .

*Case 1:  $E$  holds.* By Lemma 5.4, we may bound

$$I(X^t, \tau^t; u, \mathcal{A}^t) \leq \sum_{a \in \mathcal{A}^t} (\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)})) = O\left(\sum_{(i,j) \in X^t} \sqrt{\frac{\log(|\mathcal{A}|T)}{n_{ij}^t}}\right),$$

where  $n_{ij}^t$  is the number of times that the pair  $(i, j)$  has been matched at the start of round  $t$ . Let  $w_{i,j}^t = \frac{1}{\sqrt{n_{ij}^t}}$  be the size of the confidence set (with the log factor scaled out) for  $(i, j)$  at the start of round  $t$ .

At each timestep  $t$ , let us consider the list consisting of  $w_{i_t, j_t}^t$  for all  $(i_t, j_t) \in X^t$ . Let us now consider the overall list consisting of the concatenation of all of these lists over all rounds. Let us order this list in decreasing order to obtain a list  $\tilde{w}_1, \dots, \tilde{w}_L$  where  $L = \sum_{t=1}^T |X^t| \leq nT$ . In this notation, we observe that:

$$\sum_{t=1}^T I(X^t, \tau^t; u, \mathcal{A}) \leq \sum_{t=1}^T \sum_{a \in \mathcal{A}^t} (\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)})) = \log(|\mathcal{A}|T) \sum_{l=1}^L \tilde{w}_l.$$

We claim that  $\tilde{w}_l \leq O(\min(1, \frac{1}{\sqrt{(l/|\mathcal{A}|^2)-1}}))$ . The number of rounds that a pair of agents can have their confidence set have size at least  $\tilde{w}_l$  is upper bounded by  $1 + \frac{1}{\tilde{w}_l^2}$ . Thus, the total number of times that any confidence set can have size at least  $\tilde{w}_l$  is upper bounded by  $(|\mathcal{A}|^2)(1 + \frac{1}{\tilde{w}_l^2})$ .

Putting this together, we see that:

$$\begin{aligned} \log(|\mathcal{A}|T) \sum_{l=1}^L \tilde{w}_l &\leq O\left(\sum_{l=1}^L \min\left(1, \frac{1}{\sqrt{(l/|\mathcal{A}|^2)-1}}\right)\right) \\ &\leq O\left(\log(|\mathcal{A}|T) \sum_{l=1}^{nT} \min\left(1, \frac{1}{\sqrt{(l/|\mathcal{A}|^2)-1}}\right)\right) \\ &\leq O(|\mathcal{A}|\sqrt{nT} \log(|\mathcal{A}|T)). \end{aligned}$$

*Case 2: E does not hold.* Since each  $n_{ij}(\hat{u}_i(j) - u_i(j))$  is mean-zero and 1-subgaussian, and we have  $O(|I||\mathcal{J}|T)$  such random variables over  $T$  rounds, the probability that any of them exceeds

$$2\sqrt{\log(|I||\mathcal{J}|T/\delta)} \leq 2\sqrt{\log(|\mathcal{A}|^2T/\delta)}$$

is at most  $\delta$  by a standard tail bound for the maximum of subgaussian random variables. It follows that  $E$  fails to hold with probability at most  $|\mathcal{A}|^{-2}T^{-2}$ . In the case that  $E$  fails to hold, our regret in any given round would be at most  $4|\mathcal{A}|$  by the Lipschitz property in Proposition 4.4. (Recall that our upper confidence bound for any utility is wrong by at most 2 due to clipping each confidence interval to lie in  $[-1, 1]$ .) Thus, the expected regret from this scenario is at most

$$|\mathcal{A}|^{-2}T^{-2} \cdot 4|\mathcal{A}|T \leq 4|\mathcal{A}|^{-1}T^{-1},$$

which is negligible compared to the regret bound from when  $E$  does occur.  $\square$

## C.2 Proof of Theorem 5.2

**THEOREM 5.2.** *For preference class  $\mathcal{U}_{\text{typed}}$  (see Section 3), MATCHTYPEDUCB (defined in Section 5.3) incurs expected regret  $\mathbb{E}(R_T) = O(|C|\sqrt{nT} \log(|\mathcal{A}|T))$ , where  $n = \max_t |\mathcal{A}^t|$ .*

**PROOF OF THEOREM 5.2.** Like in the proof of Theorem 5.1, we divide into two cases, based on the event  $E$  that all of the confidence sets contain their respective true utilities at every timestep  $t \leq T$ . That is,  $u_a(a') \in C_{a, a'}$  for all pairs of agents at all  $t$ .

*Case 1: E holds.* By Lemma 5.4, we may bound

$$I(X^t, \tau^t; u, \mathcal{A}^t) \leq \sum_{a \in \mathcal{A}^t} \left( \max(C_{a, c_{\mu_{X^t}(a)}}) - \min(C_{a, c_{\mu_{X^t}(a)}}) \right) = O\left( \sum_{(i, j) \in X^t} \sqrt{\frac{\log(|\mathcal{A}|T)}{n_{c_1 c_2}^t}} \right),$$

where  $n_{c_1 c_2}^t$  is the number of times that an agent of type  $c_1$  has been matched with an agent of context  $c_2$  at the start of round  $t$ . (We define  $n_{c_1 c_2}^0 = 0$  by default.) Let  $w_{c_1, c_2}^t = \frac{1}{\sqrt{n_{c_1 c_2}^t}}$

be the size of the confidence set (with the log factor scaled out) for  $(c_1, c_2)$  at the start of round  $t$ .

At each timestep  $t$ , let us consider the list consisting of  $w_{c_{i_t}, c_{j_t}}^t$  for all  $(i_t, j_t) \in X^t$ . Let us now consider the overall list consisting of the concatenation of all of these lists over all rounds. Let us order this list in decreasing order to obtain a list  $\tilde{w}_1, \dots, \tilde{w}_L$  where  $L = \sum_{t=1}^T |X^t| \leq nT$ . In this notation, we observe that:

$$\sum_{t=1}^T I(X^t, \tau^t; u, \mathcal{A}^t) \leq \sum_{t=1}^T \sum_{a \in \mathcal{A}^t} (\max(C_{c_a, c_{\mu_{X^t}(a)}}) - \min(C_{c_a, c_{\mu_{X^t}(a)}})) = \log(|\mathcal{A}|T) \sum_{l=1}^L \tilde{w}_l.$$

We claim that  $\tilde{w}_l \leq O(\min(1, \frac{1}{\sqrt{(l/|\mathcal{C}|^2)-1}}))$ . The number of instances that a pair of contexts can have their confidence set have size at least  $\tilde{w}_l$  is upper bounded by  $2n + \frac{1}{\tilde{w}_l^2}$ . Thus, the total number of times that any confidence set can have size at least  $\tilde{w}_l$  is upper bounded by  $(|\mathcal{C}|)(2n + \frac{1}{\tilde{w}_l^2})$ .

Putting this together, we see that:

$$\begin{aligned} \log(|\mathcal{A}|T) \sum_{l=1}^L \tilde{w}_l &\leq O\left(\sum_{l=1}^L \min\left(1, \frac{1}{\sqrt{(l/|\mathcal{A}|^2)-1}}\right)\right) \\ &\leq O\left(\log(|\mathcal{A}|T) \sum_{l=1}^{nT} \min\left(1, \frac{1}{\sqrt{(l/|\mathcal{C}|^2)-1}}\right)\right) \\ &\leq O(|\mathcal{C}|\sqrt{nT} \log(|\mathcal{C}|^2T)). \end{aligned}$$

*Case 2: E does not hold.* Since each  $n_{ij}(\hat{u}_i(j) - u_i(j))$  is mean-zero and 1-subgaussian, and we have  $O(|\mathcal{I}||\mathcal{J}|T)$  such random variables over  $T$  rounds, the probability that any of them exceeds

$$2\sqrt{\log(|\mathcal{I}||\mathcal{J}|T/\delta)} \leq 2\sqrt{\log(|\mathcal{A}|^2T/\delta)}$$

is at most  $\delta$  by a standard tail bound for the maximum of subgaussian random variables. It follows that  $E$  fails to hold with probability at most  $|\mathcal{A}|^{-2}T^{-2}$ . In the case that  $E$  fails to hold, our regret in any given round would be at most  $4|\mathcal{A}|$  by the Lipschitz property in Proposition 4.4. (Recall that our upper confidence bound for any utility is wrong by at most two due to clipping each confidence interval to lie in  $[-1, 1]$ .) Thus, the expected regret from this scenario is at most

$$|\mathcal{A}|^{-2}T^{-2} \cdot 4|\mathcal{A}|T \leq 4|\mathcal{A}|^{-1}T^{-1},$$

which is negligible compared to the regret bound from when  $E$  does occur.  $\square$

### C.3 Proof of Theorem 5.3

**THEOREM 5.3.** *For preference class  $\mathcal{U}_{\text{linear}}$  (see Section 3), MATCHLINUCB (defined in Section 5.3) incurs expected regret  $\mathbb{E}(R_T) = O(d\sqrt{|\mathcal{A}|\sqrt{nT} \log(|\mathcal{A}|T)})$ , where  $n = \max_t |\mathcal{A}_t|$ .*

To prove Theorem 5.3, it suffices to (a) show that the confidence sets contain the true utilities with high probability and (b) bound the sum of the sizes of the confidence sets.

Part (a) follows from fact established in existing analysis of LinUCB in the classical linear contextual bandits setting [39].

**LEMMA C.1** ([39, PROPOSITION 2]). *Let the confidence sets be defined as above (and in MATCHLINUCB). For each  $a \in \mathcal{A}$ , it holds that:*

$$\mathbb{P}[\phi(a) \in C_{\phi(a)} \quad \forall 1 \leq t \leq T] \geq 1 - 1/(|\mathcal{A}|^3T^2).$$

LEMMA C.2. *Let the confidence sets be defined as above (and in MATCHLINUCB). For each  $a \in \mathcal{A}$  and for any  $\varepsilon > 0$ , it holds that:*

$$\sum_{t|a \in \mathcal{A}^t, \mu_{X^t}(a) \neq a} \mathbf{1} \left[ \max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)}) > \varepsilon \right] \leq O \left( \left( \frac{4\beta_T}{\varepsilon^2} + 1 \right) d \log(1/\varepsilon) \right).$$

PROOF. We follow the same argument as the proof of Proposition 3 in Russo and Van Roy [39].

We first recall the definition of  $\varepsilon$ -dependence and  $\varepsilon$ -eluder dimension: We say that an agent  $a'$  is  $\varepsilon$ -dependent on  $a'_1, \dots, a'_s$  if for all  $\phi(a), \tilde{\phi}(a) \in \mathcal{B}^d$  such that

$$\sum_{k=1}^s \langle c_{a'_k}, \tilde{\phi}(a) - \phi(a) \rangle^2 \leq \varepsilon^2,$$

we also have  $\langle c_{a'}, \tilde{\phi}(a) - \phi(a) \rangle^2 \leq \varepsilon^2$ . The  $\varepsilon$ -eluder dimension  $d_{\varepsilon\text{-eluder}}$  of  $\mathcal{B}^d$  is the maximum length of a sequence  $a'_1, \dots, a'_s$  such that no element is  $\varepsilon$ -dependent on a prefix.

Consider the subset  $S_a$  of  $\{t \mid a \in \mathcal{A}^t, \mu_{X^t}(a) \neq a\}$  such that

$$\mathbf{1} \left[ \max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)}) > \varepsilon \right].$$

Suppose for the sake of contradiction that

$$|S_a| > \left( \frac{4\beta_T}{\varepsilon^2} + 1 \right) d_{\varepsilon\text{-eluder}}.$$

Then, there exists an element  $t^*$  that is  $\varepsilon$ -dependent on  $\frac{4\beta_T}{\varepsilon^2} + 1$  disjoint subsets of  $S_a$ : One can repeatedly remove sequences  $a'_{\mu_{X^{t_1}}(a)}, \dots, a'_{\mu_{X^{t_s}}(a)}$  of maximal length such that no element is  $\varepsilon$ -dependent on a prefix; note that  $s \leq d_{\varepsilon\text{-eluder}}$  always. Let the subsets be  $S_a^{(q)}$  for  $q = 1, \dots, \frac{4\beta_T}{\varepsilon^2} + 1$ , and let  $\phi(a), \tilde{\phi}(a)$  be such that  $\langle c_{\mu_{X^{t^*}}(a)}, \tilde{\phi}(a) - \phi(a) \rangle > \varepsilon$ . The above implies that

$$\sum_{q=1}^{\frac{4\beta_T}{\varepsilon^2} + 1} \sum_{t \in S_a^{(q)}} \langle c_{\mu_{X^t}(a)}, \tilde{\phi}(a) - \phi(a) \rangle^2 > 4\beta_T$$

by the definition of  $\varepsilon$ -dependence. But this is impossible, since the left-hand side is upper bounded by

$$\sum_{t=1}^T \langle c_{\mu_{X^t}(a)}, \tilde{\phi}(a) - \phi(a) \rangle^2 \leq 4\beta_T$$

by the definition of the confidence sets. Hence, it must hold that

$$|S_a| \leq \left( \frac{4\beta_T}{\varepsilon^2} + 1 \right) d_{\varepsilon\text{-eluder}}.$$

Now, it follows from the bound on the eluder dimension for linear bandits (Proposition 6 in Reference [39]) that the bound of  $\tilde{O}((\frac{4\beta_T}{\varepsilon^2} + 1)d \log(1/\varepsilon))$  holds.  $\square$

LEMMA C.3. *Let the confidence sets be defined as above (and in MATCHLINUCB). For any  $a \in \mathcal{A}$ , it holds that:*

$$\sum_{t|a \in \mathcal{A}^t, \mu_{X^t}(a) \neq a} (\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)})) \leq O(d(\log(T|\mathcal{A}|))\sqrt{T_a}),$$

where  $T_a$  is the number of times that agent  $a$  has been matched.

PROOF. Let us consider the set of confidence set sizes  $(\max(C_{a, \mu_{X^t}(a)} - \min(C_{a, \mu_{X^t}(a)})))$  for  $t$  such that  $a \in \mathcal{A}^t, \mu_{X^t}$ . Let us sort these confidence set sizes in decreasing order and label them  $w_1 \geq \dots \geq w_{T_a}$ . Restating Lemma C.2, we see that

$$\sum_{t=1}^{T_a} w_t \mathbf{1}[w_t > \varepsilon] \leq O\left(\left(\frac{4\beta_T}{\varepsilon^2} + 1\right) d \log(1/\varepsilon)\right). \quad (10)$$

for all  $\varepsilon > 0$ .

We see that:

$$\begin{aligned} \sum_{t|a \in \mathcal{A}^t, \mu_{X^t}(a) \neq a} (\max(C_{a, \mu_{X^t}(a)} - \min(C_{a, \mu_{X^t}(a)}))) &= \sum_{t=1}^{T_a} w_t \\ &\leq \sum_{t=1}^{T_a} w_t \mathbf{1}[w_t > 1/T_a^2] + \sum_{t=1}^{T_a} w_t \mathbf{1}[w_t \leq 1/T_a^2] \\ &\leq \frac{1}{T_a} + \sum_{t=1}^{T_a} w_t \mathbf{1}[w_t > 1/T_a^2]. \end{aligned}$$

We claim that  $w_i \leq 2$  if  $i \geq d \log(T_a)$  and  $w_i \leq \min(2, \frac{4\beta_T(d \log T_a)}{i - d \log T_a})$  if  $i < d \log T_a$ . The first part follows from the fact that we truncate the confidence sets to be within  $[-1, 1]$ . It thus suffices to show that  $w_i \leq \frac{4\beta_T(d \log T_a)}{i - d \log T_a}$  for  $t \leq d \log T$ . If  $w_i \geq \varepsilon > 1/T_a^2$ , then we see that  $\sum_{t=1}^{T_a} \mathbf{1}[w_t > \varepsilon] \geq i$ , which means by (10) that  $i \leq O((\frac{4\beta_T}{\varepsilon^2} + 1)d \log(1/\varepsilon)) \leq O((\frac{4\beta_T}{\varepsilon^2} + 1)d \log(T_a))$ , which means that  $\varepsilon \leq \frac{4\beta_T(d \log T_a)}{i - d \log T_a}$ . This proves the desired statement.

Now, we can plug this into the above expression to obtain:

$$\begin{aligned} &\sum_{t|a \in \mathcal{A}^t, \mu_{X^t}(a) \neq a} (\max(C_{a, \mu_{X^t}(a)} - \min(C_{a, \mu_{X^t}(a)}))) \\ &\leq \frac{1}{T_a} + \sum_{t=1}^{T_a} w_t \mathbf{1}[w_t > 1/T_a^2] \\ &\leq \frac{1}{T_a} + 2d \log(T_a) + \sum_{i > d \log T_a} \min\left(2, \frac{4\beta_T(d \log T_a)}{i - d \log T_a}\right) \\ &\leq \frac{1}{T_a} + 2d \log(T_a) + 2\sqrt{d \log T_a \beta_T} \int_{t=0}^{T_a} t^{-1/2} dt \\ &= \frac{1}{T_a} + 2d \log(T_a) + 4\sqrt{d T_a \log T_a \beta_T}. \end{aligned}$$

We now use the fact that:

$$\beta_T = O\left(d \log T + \frac{1}{T} \sqrt{\log(T^2 |A|)}\right).$$

Plugging this into the above expression, we obtain the desired result.  $\square$

We are now ready to prove Theorem 5.3.

PROOF OF THEOREM 5.3. Like in the proof of Theorem 5.1, we divide into two cases, based on the event  $E$  that all of the confidence sets contain their respective true utilities at every timestep  $t \leq T$ . That is,  $u_{c_1}(c_2) \in C_{c_1, c_2}$  for all  $c_1, c_2 \in C$  at all  $t$ .



*Case 1:  $E$  holds.* By Lemma 5.4, we know that the cumulative regret is upper bounded by

$$\begin{aligned} R_T &\leq \sum_{t=1}^T \sum_{a \in \mathcal{A}^t} (\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)})) \\ &= \sum_{a \in \mathcal{A}} \sum_{t | a \in \mathcal{A}^t, \mu_{X^t}(a) \neq a} (\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)})) \\ &\leq \sum_{a \in \mathcal{A}} O(d \log(T|\mathcal{A}|) \sqrt{T_a}), \end{aligned}$$

where the last inequality applies Lemma C.3 to the inner summand. We see that  $\sum_{a \in \mathcal{A}} T_a = \sum_t |\mathcal{A}_t| \leq nT$  by definition, since at most  $n$  agents show up at every round. Let us now observe that:

$$\sum_{a \in \mathcal{A}} \sqrt{T_a} \leq \sqrt{|\mathcal{A}|} \sqrt{\sum_{a \in \mathcal{A}} T_a} \leq \sqrt{|\mathcal{A}| nT},$$

as desired.

*Case 2:  $E$  does not hold.* From Lemma C.1, it follows that:

$$\mathbb{P}[\phi(a) \in C_{\phi(a)} \quad \forall 1 \leq t \leq T] \geq 1 - 1/(|\mathcal{A}|^3 T^2).$$

Union bounding, we see that

$$\mathbb{P}[\phi(a) \in C_{\phi(a)} \quad \forall 1 \leq t \leq T \forall a \in \mathcal{A}] \geq 1 - 1/(|\mathcal{A}|^2 T^2).$$

By the definition of the confidence sets for the utilities, we see that:

$$\mathbb{P}[u(a, a') \in C_{a, a'} \quad \forall 1 \leq t \leq T, \forall a, a' \in \mathcal{A}] \geq 1/(|\mathcal{A}|^2 T^2). \quad (11)$$

Thus, the probability that event  $E$  does not hold is at most  $|\mathcal{A}|^{-2} T^{-2}$ . In the case that  $E$  fails to hold, our regret in any given round would be at most  $4|\mathcal{A}|$  by the Lipschitz property in Proposition 4.4. Thus, the expected regret is at most  $4|\mathcal{A}|^{-1} T^{-1}$ , which is negligible compared to the regret bound from when  $E$  does occur.  $\square$

#### C.4 Proof of Lemma 5.5

**LEMMA 5.5.** *For any algorithm that learns a stable matching with respect to unstructured preferences, there exists an instance on which it has expected regret  $\tilde{\Omega}(|\mathcal{A}|^{3/2} \sqrt{T})$  (where regret is given by Subset Instability).*

**PROOF OF LEMMA 5.5.** Recall that, by Proposition 4.4, the problem of learning a maximum weight matching with respect to utility difference is no harder than that of learning a stable matching with respect to Subset Instability. In the remainder of our proof, we reduce a standard “hard instance” for stochastic multi-armed bandits to our setting of learning a maximum weight matching.

*Step 1: Constructing the hard instance for stochastic MAB.* Consider the following family of stochastic multi-armed bandits instances: For a fixed  $K$ , let  $\mathcal{I}_\alpha$  for  $\alpha \in \{1, \dots, K\}$  denote the stochastic multi-armed bandits problem where all arms have 0-1 rewards, and the  $k$ -th arm has mean reward  $\frac{1}{2} + \rho$  if  $k = \alpha$  and  $\frac{1}{2}$  otherwise, where  $\rho > 0$  will be set later. A classical lower bound for stochastic multi-armed bandits is the following:

**LEMMA C.4 ([6]).** *The expected regret of any stochastic multi-armed bandit algorithm on an instance  $\mathcal{I}_\alpha$  for  $\alpha$  selected uniformly at random from  $\{1, \dots, K\}$  is  $\Omega(\sqrt{KT})$ .*

*Step 2: Constructing a (random) instance for the maximum weight matching problem.* We will reduce solving the above distribution over stochastic multi-armed bandits problems to a distribution over instances of learning a maximum weight matching. Let us now construct this random instance of the maximum weight matching problem. Let  $|\mathcal{I}| = K$  and  $|\mathcal{J}| = 10K \log(KT)$ . Specifically, we sample inputs for learning a maximum weight matching as follows: For each man  $i \in \mathcal{I}$ , select  $\alpha_i \in \{1, \dots, K\}$  uniformly at random, and define  $u_i(j)$  to be  $\frac{1}{2} + \rho$  if  $\lfloor (j-1)/\log K \rfloor = \alpha_i$  and  $\frac{1}{2}$  otherwise. Furthermore, let  $u_j(i) = 0$  for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$ . Finally, suppose observations are always in  $\{0, 1\}$  (but are unbiased).

The key property of the above setup that we will exploit for our reduction is the fact that, due to the imbalance in the market, the maximum weight matching for these utilities has with high probability each  $i$  matched with some  $j$  whom they value at  $\frac{1}{2} + \rho$ . Indeed, by a union bound, the probability that more than  $10 \log(KT)$  different  $i$  have the same  $\alpha_i$  is at most

$$K \cdot \binom{K}{10 \log(KT)} K^{-10 \log(KT)} = O(K^{-4} T^{-4}).$$

Thus, with probability  $1 - O(K^{-4} T^{-4})$ , this event holds. The case where this event does not hold contributes negligibly to regret, so we do not consider it further.

*Step 3: Establishing the reduction.* Now, suppose for the sake of contradiction that some algorithm could solve our random instance of learning a maximum weight matching problem with expected regret  $o(K^{3/2} \sqrt{T})$ . We can obtain a stochastic multi-armed bandits that solves the instances in Lemma C.4 as follows: Choose a random  $i^* \in \mathcal{I}$  and set  $\alpha_{i^*} = \alpha$ . Simulate the remaining  $i$  by choosing  $\alpha_i$  for all  $i \neq i^*$  uniformly at random. Run the algorithm on this instance of learning a maximum weight matching, “forwarding” arm pulls to the true instance when matching  $i^*$ .

To analyze the regret of this algorithm when faced with the distribution from Lemma C.4, we first note that with high probability, all the agents  $i \in \mathcal{I}$  can simultaneously be matched to a set of  $j \in \mathcal{J}$  such that each  $i$  is matched to some  $j$  whom they value at  $\frac{1}{2} + \rho$ . Then, the regret of any matching is  $\rho$  times the number of  $i \in \mathcal{I}$  who are not matched to a  $j$  whom they value at  $\frac{1}{2} + \rho$ . Thus, we can define the cumulative regret for an agent  $i \in \mathcal{I}$  as  $\rho$  times the number of rounds they were not matched to someone whom they value at  $\frac{1}{2} + \rho$ . For  $i^*$ , this regret is just the regret for the distribution from Lemma C.4. Since  $i^*$  was chosen uniformly at random, their expected cumulative regret is at most

$$\frac{1}{K} \cdot o(K^{3/2} \sqrt{T}) = o(\sqrt{KT}),$$

in violation of Lemma C.4.

*Step 4: Concluding the lower bound.* This contradiction implies that no algorithm can hope to obtain  $o(K^{3/2} \sqrt{T})$  expected regret on this distribution over instances of learning a maximum weight matching. Since there are  $O(K \log(KT)) = \tilde{O}(K)$  agents in the market total, the desired lower bound follows.  $\square$

## D PROOF OF THEOREM 6.1

**THEOREM 6.1 (INSTANCE-DEPENDENT REGRET).** *Suppose that  $\mathcal{A}_t = \mathcal{A}$  for all  $t$ . Let  $u \in \mathcal{U}_{\text{unstructured}}$  be any utility function and put*

$$\Delta := \inf_{X \neq X^*} \left\{ \sum_{a \in \mathcal{A}} u_a(\mu_{X^*}(a)) - \sum_{a \in \mathcal{A}} u_a(\mu_X(a)) \right\}.$$

*Then, MATCHUCB' incurs expected regret  $\mathbb{E}(R_T) = O(|\mathcal{A}|^5 \cdot \log(|\mathcal{A}|T)/\Delta^2)$ .*

### D.1 MATCHUCB'

MATCHUCB' is the same as MATCHUCB, except we call COMPUTEMATCH' instead of COMPUTEMATCH. The idea behind COMPUTEMATCH' is that we compute an optimal primal-dual solution for both the original confidence sets  $C$  as well as expanded confidence sets  $C'$ , which we define to be twice the width of the original confidence sets. More formally, we define

$$C'_{a,a'} := \left[ \min(C_{a,a'}) - \frac{\max(C_{a,a'}) - \min(C_{a,a'})}{2}, \max(C_{a,a'}) + \frac{\max(C_{a,a'}) - \min(C_{a,a'})}{2} \right].$$

We will adaptively explore (following UCB) according to both  $C$  and  $C'$ . Doing extra exploration according to the more pessimistic confidence sets  $C'$  is necessary for us to be able to find “robust” dual solutions for setting transfers.

We define  $(X^*, p^*)$ , which will be an optimal primal-dual solution for the upper confidence bounds of  $C$  as follows: Let  $X^*$  be a maximum weight matching with respect to  $u^{\text{UCB}}$ . We next compute the gap

$$\Delta^{\text{UCB}} = \min_{X \neq X^*} \left\{ \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)) - \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_X(a)) \right\}$$

with respect to  $u^{\text{UCB}}$ . We can compute this gap by computing the maximum weight matching and the second-best matching with respect to  $u^{\text{UCB}}$ .<sup>15</sup> Next, define utility functions  $u'_a$  such that

$$u'_a(a') = \begin{cases} u_a^{\text{UCB}}(a') - \frac{\Delta^{\text{UCB}}}{|\mathcal{A}|} & \text{if } \mu_{X^*}(a) = a' \text{ and } a \neq a' \\ u_a^{\text{UCB}}(a') & \text{otherwise} \end{cases}$$

for all  $a \in \mathcal{A}$ . (We show in Lemma D.1 that  $X^*$  is still a maximum weight matching for  $u'$ .) Now, compute an optimal dual solution  $p'$  for utility function  $u'$ . To get  $p^*$ , we add  $\Delta^{\text{UCB}}/|\mathcal{A}|$  to  $p'_a$  for each matched agent  $a$  in  $X^*$ . (See Lemma D.2 for a proof that  $(X^*, p^*)$  is an optimal primal-dual pair with respect to  $u^{\text{UCB}}$ .)

Finally, let  $(X^{*,2}, p^{*,2})$  be any optimal primal-dual pair for the utility function  $u^{\text{UCB},2}$  given by the upper confidence bounds  $\max(C'_{a,a'})$  of  $C'$ .

With this setup, we define COMPUTEMATCH' as follows: If  $X^* \neq X^{*,2}$ , then return  $(X^{*,2}, \tau^{*,2})$ , where  $\tau^{*,2}$  is given by  $\tau_a^{*,2} = p_a^{*,2} - u_a^{\text{UCB},2}(\mu_{X^{*,2}}(a))$  if  $a$  is matched and  $\tau_a^{*,2} = 0$  if  $a$  is unmatched. Otherwise, return  $(X^*, \tau^*)$ , where  $\tau^*$  is given by  $\tau_a^* = p_a^* - u_a^{\text{UCB}}(\mu_X(a))$  if  $a$  is matched and  $\tau_a^* = 0$  if  $a$  is unmatched.

### D.2 Proof of Theorem 6.1

We first verify (as claimed above) that  $X^*$  is a maximum weight matching with respect to  $u'$ .

LEMMA D.1. *Matching  $X^*$  is a maximum weight matching with respect to  $u'$ .*

PROOF. Consider any matching  $X \neq X^*$ . Since

$$\sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_X(a)) \leq -\Delta^{\text{UCB}} + \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a))$$

by the definition of  $\Delta^{\text{UCB}}$ , we have

$$\sum_{a \in \mathcal{A}} u'_a(\mu_X(a)) \leq \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_X(a)) \leq \sum_{a \in \mathcal{A}} (u_a^{\text{UCB}}(\mu_{X^*}(a)) - \frac{\Delta^{\text{UCB}}}{|\mathcal{A}|}) \leq \sum_{a \in \mathcal{A}} u'_a(\mu_{X^*}(a)). \quad \square$$

<sup>15</sup>See Chegireddy and Hamacher [14] for efficient algorithms to compute the second-best matching.

We now prove the main lemma for this analysis, restated below. Lemma 6.2 shows that if the confidence sets are small enough, then the selected matching will be stable with respect to the true utilities.

LEMMA 6.2. *Suppose COMPUTEMATCH' is run on a collection  $\mathcal{C}$  of confidence sets  $C_{i,j}$  and  $C_{j,i}$  over the agent utilities that satisfy*

$$\max(C_{i,j}) - \min(C_{i,j}) \leq 0.05 \frac{\Delta}{|\mathcal{A}|} \quad \text{and} \quad \max(C_{j,i}) - \min(C_{j,i}) \leq 0.05 \frac{\Delta}{|\mathcal{A}|}$$

*for all  $(i,j)$  in the matching returned by COMPUTEMATCH'. Suppose also that the confidence sets  $\mathcal{C}$  contain the true utilities for all pairs of agents. Then, the market outcome returned by COMPUTEMATCH' is stable with respect to the true utilities  $u$ .*

PROOF OF LEMMA 6.2. The proof proceeds in five steps, which we now outline. We first show the matching returned by COMPUTEMATCH' is the maximum weight matching  $X^{\text{opt}}$  with respect to  $u$ . We next show that  $X^*$  as defined in COMPUTEMATCH' also equals  $X^{\text{opt}}$ . These facts let us conclude that COMPUTEMATCH' returns  $(X^*, \tau^*)$ . We then show  $\Delta^{\text{UCB}}$  is at least  $0.1\Delta$ . We then show that  $(X^*, \tau^*)$  is stable with respect to  $u'$ . We finish by showing that this implies  $(X^*, \tau^*)$  is stable with respect to  $u$ .

Throughout the proof, we will use the following observation about the expanded confidence sets:

$$\max(C'_{i,j}) - \min(C'_{i,j}) \leq 0.1 \frac{\Delta}{|\mathcal{A}|} \quad \text{and} \quad \max(C'_{j,i}) - \min(C'_{j,i}) \leq 0.1 \frac{\Delta}{|\mathcal{A}|} \quad (12)$$

for all  $(i,j)$  in the matching returned by COMPUTEMATCHING'. This follows from the assumptions in the lemma statement.

*Proving COMPUTEMATCH' returns  $X^{\text{opt}}$  as the matching.* COMPUTEMATCH' by definition returns  $X^{*,2}$  always, so it suffices to show that  $X^{*,2} = X^{\text{opt}}$ . Note that  $X^{*,2}$  is a maximum weight matching with respect to  $u^{\text{UCB},2}$ . This means that

$$\begin{aligned} \sum_{a \in \mathcal{A}} u_a(\mu_{X^{*,2}}(a)) &\geq - \sum_{a \in \mathcal{A}} (\max(C'_{a,\mu_{X^{*,2}}(a)}) - \min(C'_{a,\mu_{X^{*,2}}(a)})) + \sum_{a \in \mathcal{A}} u_a^{\text{UCB},2}(\mu_{X^{*,2}}(a)) \\ &\geq -0.1\Delta + \sum_{a \in \mathcal{A}} u_a^{\text{UCB},2}(\mu_{X^{*,2}}(a)) \\ &\geq -0.1\Delta + \sum_{a \in \mathcal{A}} u_a^{\text{UCB},2}(\mu_{X^{\text{opt}}}(a)) \\ &\geq -0.1\Delta + \sum_{a \in \mathcal{A}} u_a(\mu_{X^{\text{opt}}}(a)). \end{aligned}$$

By the definition of the gap  $\Delta$ , we conclude that  $X^{*,2} = X^{\text{opt}}$ .

*Proving  $X^* = X^{\text{opt}}$ .* Suppose for sake of contradiction that  $X^* \neq X^{\text{opt}}$ . Then,

$$\sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)) \geq \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^{\text{opt}}}(a)) \geq \sum_{a \in \mathcal{A}} u_a(\mu_{X^{\text{opt}}}(a)),$$

since  $X^*$  is a maximum weight matching with respect to  $u^{\text{UCB}}$ . Moreover, by the definition of the gap, we know that  $\sum_{a \in \mathcal{A}} u_a(\mu_{X^*}(a)) \leq \sum_{a \in \mathcal{A}} u_a(\mu_{X^{\text{opt}}}(a)) - \Delta$ . Putting this all together, we see

that

$$\begin{aligned} \sum_{a \in \mathcal{A}} (\max(C_{a, \mu_{X^*}(a)}) - \min(C_{a, \mu_{X^*}(a)})) &\geq \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)) - \sum_{a \in \mathcal{A}} u_a(\mu_{X^*}(a)) \\ &\geq \Delta. \end{aligned}$$

We now use this to lower bound the utility of  $X^*$  on  $u^{\text{UCB}, 2}$ . By the definition of the confidence sets, we see that

$$\begin{aligned} \sum_{a \in \mathcal{A}} u_a^{\text{UCB}, 2}(\mu_{X^*}(a)) &\geq \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)) + \frac{1}{2} \sum_{a \in \mathcal{A}} (\max(C_{a, \mu_{X^*}(a)}) - \min(C_{a, \mu_{X^*}(a)})) \\ &\geq \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)) + 0.5\Delta. \end{aligned}$$

However,  $X^{\text{opt}}$  only achieves a utility of

$$\begin{aligned} \sum_{a \in \mathcal{A}} u_a^{\text{UCB}, 2}(\mu_{X^{\text{opt}}}(a)) &\leq \sum_{a \in \mathcal{A}} u_a(\mu_{X^{\text{opt}}}(a)) + \sum_{a \in \mathcal{A}} (\max(C'_{a, \mu_{X^{\text{opt}}}(a)}) - \min(C'_{a, \mu_{X^{\text{opt}}}(a)})) \\ &\leq \sum_{a \in \mathcal{A}} u_a(\mu_{X^{\text{opt}}}(a)) + 0.1\Delta. \end{aligned}$$

But this contradicts the fact (from above) that  $X^{\text{opt}} = X^{*, 2}$  is a maximum weight matching with respect to  $u^{\text{UCB}, 2}$ . Therefore, it must be that  $X^* = X^{\text{opt}}$ .

Putting the above two arguments together, we conclude  $\text{COMPUTEMATCH}'$  returns  $(X^*, \tau^*)$  in this case.

*Bounding the gap  $\Delta^{\text{UCB}}$ .* We next show that  $\Delta^{\text{UCB}} \geq 0.1\Delta$ . We proceed by assuming

$$\sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_X(a)) \geq -0.1\Delta + \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)) \quad (13)$$

for some  $X \neq X^*$  and deriving a contradiction.

We first show that (13) implies a lower bound on

$$S = \sum_{a \in \mathcal{A}} (\max(C_{a, \mu_X(a)}) - \min(C_{a, \mu_X(a)}))$$

in terms of  $\Delta$ . Because the confidence sets contain the true utilities and  $u_a^{\text{UCB}}$  upper bounds  $u_a$  pointwise, (13) implies

$$S + \sum_{a \in \mathcal{A}} u_a(\mu_X(a)) \geq \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_X(a)) \geq -0.1\Delta + \sum_{a \in \mathcal{A}} u_a(\mu_{X^*}(a)).$$

Applying the definition of  $\Delta$ , we obtain the lower bound

$$S \geq -0.1\Delta + \sum_{a \in \mathcal{A}} u_a(\mu_{X^*}(a)) - \sum_{a \in \mathcal{A}} u_a(\mu_X(a)) \geq (1 - 0.1)\Delta.$$

Now, we apply the fact that  $X^* = X^{*, 2} = X^{\text{opt}}$ . We establish the following contradiction:

$$\begin{aligned} 0.1\Delta + \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)) &\geq 0.1\Delta + \sum_{a \in \mathcal{A}} u_a(\mu_{X^*}(a)) \\ &= \sum_{a \in \mathcal{A}} (u_a(\mu_{X^*}(a)) + 0.1\Delta/|\mathcal{A}|) \\ &\stackrel{(i)}{\geq} \sum_{a \in \mathcal{A}} u_a^{\text{UCB}, 2}(\mu_{X^*}(a)) \end{aligned}$$

$$\begin{aligned}
&\stackrel{(ii)}{\geq} \sum_{a \in \mathcal{A}} u_a^{\text{UCB},2}(\mu_X(a)) \\
&\stackrel{(iii)}{\geq} \frac{S}{2} + \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_X(a)) \\
&\stackrel{(iv)}{\geq} \left(\frac{1}{2}(1 - 0.1)\right) \Delta + \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_X(a)) \\
&\stackrel{(v)}{\geq} \left(\frac{1}{2}(1 - 0.1) - 0.1\right) \Delta + \sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)).
\end{aligned}$$

Here, (i) comes from (12) in the lemma statement; (ii) holds because  $X^* = X^{*,2}$  is a maximum weight matching with respect to  $u^{\text{UCB},2}$ ; (iii) is by the definition of  $u^{\text{UCB},2}$ ; (iv) follows from our lower bound on  $S$ ; and (v) follows from (13).

*Proving that  $(X^*, \tau^*)$  is stable with respect to  $u'$ .* By Lemma D.1,  $(X^*, p')$  is an optimal primal-dual pair with respect to  $u'$ . Now, it suffices to show that the primal-dual solution corresponds to the market outcome  $(X^*, \tau^*)$  for  $u'$ . To see this, notice that  $p'_a = 0$  for unmatched agents and

$$p'_a = p_a^* - \frac{\Delta^{\text{UCB}}}{2|\mathcal{A}|} = \tau_a^* + u'_a(\mu_{X^*}(a))$$

for matched agents.

*Proving that  $(X^*, \tau^*)$  is stable with respect to  $u$ .* We show the stability  $(X^*, \tau^*)$  with respect to  $u$  by checking that individual rationality holds and that there are no blocking pairs.

The main fact that we will use is that

$$u_a(\mu_{X^*}(a)) \geq u'_a(\mu_{X^*}(a)).$$

To prove this, we split into two cases: (i) agent  $a$  is matched in  $X^*$  (i.e.,  $\mu_{X^*}(a) \neq a$ ), and (ii) agent  $a$  is not matched by  $X^*$ . For (i), if  $a$  is matched by  $X^*$ , then

$$u_a(\mu_{X^*}(a)) \geq u_a^{\text{UCB}}(\mu_{X^*}(a)) - 0.1 \frac{\Delta}{|\mathcal{A}|} \geq u_a^{\text{UCB}}(\mu_{X^*}(a)) - \frac{\Delta^{\text{UCB}}}{|\mathcal{A}|} = u'_a(\mu_{X^*}(a)).$$

For (ii), if  $a$  is not matched by  $X^*$ , then  $u_a(\mu_{X^*}(a)) \geq u'_a(\mu_{X^*}(a))$ , because both sides are 0.

For individual rationality, we thus have

$$u_a(\mu_{X^*}(a)) + \tau_a^* \geq u'_a(\mu_{X^*}(a)) + \tau_a^* \geq 0,$$

where the second inequality comes from the individual rationality of  $(X^*, \tau^*)$  with respect to  $u'$ .

Let us next show that there are no blocking pairs. If  $(i, j) \in X^*$ , then we see that:

$$u_i(\mu_{X^*}(i)) + \tau_i^* + u_j(\mu_{X^*}(j)) + \tau_j^* = u_i(\mu_{X^*}(i)) + u_j(\mu_{X^*}(j)),$$

as desired. Next, consider any pair  $(i, j) \notin X^*$ . Then,

$$u_i(j) + u_j(i) \leq u_i^{\text{UCB}}(j) + u_j^{\text{UCB}}(i) = u'_i(j) + u'_j(i).$$

It follows that

$$\begin{aligned}
u_i(\mu_{X^*}(i)) + \tau_i^* + u_j(\mu_{X^*}(j)) + \tau_j^* &\geq u'_i(\mu_{X^*}(i)) + \tau_i^* + u_j(\mu_{X^*}(j)) + \tau_j^* \\
&\geq u'_i(j) + u'_j(i) \\
&\geq u_i(j) + u_j(i),
\end{aligned}$$



where the second inequality comes from the fact that  $(X^*, \tau^*)$  has no blocking pairs with respect to  $u'$ .

This completes our proof that  $(X^*, \tau^*)$  is stable with respect to  $u$ .  $\square$

Now, we are ready to prove Theorem 6.1.

**PROOF OF THEOREM 6.1.** As in the proof of Theorem 5.1, the starting point for our proof is the typical approach in multi-armed bandits and combinatorial bandits [15, 22, 30] of bounding regret in terms of the sizes of the confidence interval of the chosen arms. Our approach does not quite compose cleanly with these proofs, since we need to handle the transfers in addition to the matching.

We divide in two cases, based on the event  $E$  that all of the confidence sets contain their respective true utilities at every timestep  $t \leq T$ . That is,  $u_i(j) \in C_{i,j}$  and  $u_j(i) \in C_{j,i}$  for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$  at all  $t$ .

*Case 1:  $E$  holds.* Let  $n_{ij}^t$  be the number of times that the pair  $(i, j)$  has been matched by round  $t$ . For each pair  $(i, j)$ , we maintain a “blame” counter  $b_{ij}^t$ . We will ultimately bound the total number of timesteps where the algorithm chooses a matching that is not stable by  $\sum_{(i,j)} b_{ij}^T$ .

We increment the blame counters as follows. First, suppose that

$$\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)}) \leq 0.1 \frac{\Delta}{|\mathcal{A}|}$$

for every matched agent  $a \in \mathcal{A}$ . By Lemma 6.2 and since the event  $E$  holds, we know the chosen matching is stable and thus incurs 0 regret. We do not increment any of the blame counters in this case. Now, suppose that

$$\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)}) > 0.1 \frac{\Delta}{|\mathcal{A}|}$$

for some matched agent  $a$ . We increment the counter of the least-blamed pair  $(i, j) \in X^t$ .

We now bound the blame counter  $b_{ij}^T$ . We use the fact that the blame counter is only incremented when the corresponding confidence set is sufficiently large, and that a new sample of the utilities is received whenever the blame counter is incremented. This means that:

$$b_{ij}^T = O\left(\frac{|\mathcal{A}|^2 \log(|\mathcal{A}|T)}{\Delta^2}\right).$$

The maximum regret incurred by any matching is at most  $12|\mathcal{A}|$ , which means that the regret incurred by this case is at most:

$$12|\mathcal{A}| \sum_{(i,j)} b_{ij}^T \leq 12|\mathcal{A}| \sum_{(i,j)} O\left(\frac{|\mathcal{A}|^2 \log(|\mathcal{A}|T)}{\Delta^2}\right) = O\left(\frac{|\mathcal{A}|^5 \log(|\mathcal{A}|T)}{\Delta^2}\right).$$

*Case 2:  $E$  does not hold.* Since each  $n_{ij}(\hat{u}_i(j) - u_i(j))$  is mean-zero and 1-subgaussian and we have  $O(|\mathcal{I}||\mathcal{J}|T)$  such random variables over  $T$  rounds, the probability that any of them exceeds

$$2\sqrt{\log(|\mathcal{I}||\mathcal{J}|T/\delta)} \leq 2\sqrt{\log(|\mathcal{A}|^2 T/\delta)}$$

is at most  $\delta$  by a standard tail bound for the maximum of subgaussian random variables. It follows that  $E$  fails to hold with probability at most  $|\mathcal{A}|^{-2}T^{-2}$ . In the case that  $E$  fails to hold, our regret in any given round would be at most  $12|\mathcal{A}|$  by the Lipschitz property in Proposition 4.4. (Recall that our upper confidence bound is off by at most 6 due to clipping the confidence interval to lie in

$[-1, 1]$ , so the expanded confidence sets also necessarily lie in  $[-3, 3]$ .) Thus, the expected regret from this scenario is at most

$$|\mathcal{A}|^{-2}T^{-2} \cdot 12|\mathcal{A}|T \leq 12|\mathcal{A}|^{-1}T^{-1},$$

which is negligible compared to the regret bound from when  $E$  does occur.  $\square$

### D.3 Instance-independent Regret Bounds for MATCHUCB'

To establish *instance-independent* regret bounds for MATCHUCB', we show that  $(X^*, p^*)$  is indeed optimal with respect to  $u^{\text{UCB}}$ ; the remainder then follows the same argument as Theorem 5.1.

LEMMA D.2. *The pair  $(X^*, p^*)$  is an optimal primal-dual pair with respect to  $u^{\text{UCB}}$ .*

PROOF. It suffices to verify feasibility and, by weak duality, check that  $X^*$  and  $p^*$  achieve the same objective value. It is clear that  $X^*$  is primal feasible. For dual feasibility, if  $(i, j) \notin X^*$ , then

$$p_i^* + p_j^* \geq p'_i + p'_j \geq u'_i(j) + u'_j(i) = u_i^{\text{UCB}}(j) + u_j^{\text{UCB}}(i);$$

and if  $(i, j) \in X^*$ , then

$$p_i^* + p_j^* = p'_i + p'_j + 2 \frac{\Delta^{\text{UCB}}}{|\mathcal{A}|} \geq u'_i(j) + u'_j(i) + 2 \frac{\Delta^{\text{UCB}}}{|\mathcal{A}|} = u_i^{\text{UCB}}(j) + u_j^{\text{UCB}}(i).$$

Finally, we check that they achieve the same objective value with respect to  $u^{\text{UCB}}$ . By Lemma D.1 and strong duality,  $X^*$  achieves the same objective value as  $p'$  with respect to  $u'$ . Hence,

$$\sum_{a \in \mathcal{A}} u_a^{\text{UCB}}(\mu_{X^*}(a)) = 2|X^*| \frac{\Delta^{\text{UCB}}}{|\mathcal{A}|} + \sum_{a \in \mathcal{A}} u'_a(\mu_{X^*}(a)) = 2|X^*| \frac{\Delta^{\text{UCB}}}{|\mathcal{A}|} + \sum_{a \in \mathcal{A}} p'_a = \sum_{a \in \mathcal{A}} p_a^*. \quad \square$$

## E PROOFS FOR SECTION 6.2

THEOREM 6.3. *For preference class  $\mathcal{U}_{\text{unstructured}}$  (see Section 3), there exists an algorithm giving the platform*

$$\varepsilon T \sum_{t=1}^T |\mathcal{A}_t| - O\left(|\mathcal{A}| \sqrt{nT} \sqrt{\log(|\mathcal{A}||T|)}\right)$$

*revenue in the presence of search frictions while maintaining stability with high probability.*

PROOF OF THEOREM 6.3. The algorithm is defined as follows: We set confidence sets according to MATCHUCB and run essentially that algorithm, but with a modified COMPUTEMATCH. Instead of COMPUTEMATCH, we use the following algorithm. The platform first computes a matching with transfers  $(X^*, \tau^*)$  according to the UCB estimates  $u^{\text{UCB}}$ , like before. Then, the platform chooses  $X^*$  to be the selected matching and sets the transfers according to:

$$\tau_a = \tau_a^* - \varepsilon + \max(C_{a, \mu_X(a)}) - \min(C_{a, \mu_X(a)}).$$

This choice of transfers has a clean economic intuition: Agents should be compensated based on the platform's uncertainty about their utilities with  $\varepsilon$  of their transfer shaved off as revenue for the platform.

First, we show that if the confidence sets contain the true utilities, then  $(X^*, \tau)$  is  $\varepsilon$ -stable. It suffices to show that  $(X^*, \tau')$  where:

$$\tau'_a = \tau_a^* + \max(C_{a, \mu_X(a)}) - \min(C_{a, \mu_X(a)})$$

is stable. First, we see that

$$u_a(\mu_{X^{\text{UCB}}}(a)) + \tau'_a = u_a^{\text{UCB}}(\mu_{X^{\text{UCB}}}(a)) + \tau_a^* \geq 0,$$

since  $(X, \tau^*)$  is stable with respect to  $u^{\text{UCB}}$ . Furthermore, we see that:

$$\begin{aligned} (u_i(\mu_X(i)) + \tau'_i) + (u_j(\mu_X(j)) + \tau'_j) &\geq (u_i^{\text{UCB}}(\mu_X(i)) + \tau_i^*) + (u_j^{\text{UCB}}(\mu_X(j)) + \tau_j^*) \\ &\geq u_i^{\text{UCB}}(j) + u_j^{\text{UCB}}(i) \\ &\geq u_i(j) + u_j(i), \end{aligned}$$

where the second-to-last line follows from the fact that  $(X, \tau^*)$  is stable with respect to  $u^{\text{UCB}}$ .

We first show that  $s$  is a feasible solution to  $(\dagger)$ :

$$\begin{aligned} &\min(u_i(j) - u_i(\mu_{X^{\text{UCB}}}(i)) - s_i, \tilde{u}_j(i) - u_j(\mu_{X^{\text{UCB}}}(j)) - s_j) \\ &= \min(u_i(j) - u_i^{\text{UCB}}(\mu_{X^{\text{UCB}}}(i)), \tilde{u}_j(i) - u_j^{\text{UCB}}(\mu_{X^{\text{UCB}}}(j))) \\ &\leq \min(u_i^{\text{UCB}}(j) - u_i^{\text{UCB}}(\mu_{X^{\text{UCB}}}(i)), u_j^{\text{UCB}}(i) - u_j^{\text{UCB}}(\mu_{X^{\text{UCB}}}(j))) \\ &\leq 0, \end{aligned}$$

where the last step uses the fact that  $\mu_{X^{\text{UCB}}}$  is stable with respect to  $u^{\text{UCB}}$  by definition. Moreover, we see that

$$u_a(\mu_{X^{\text{UCB}}}(a)) + s_a = u_a^{\text{UCB}}(\mu_{X^{\text{UCB}}}(a)) \geq 0,$$

where the last inequality uses that  $\mu_{X^{\text{UCB}}}$  is stable with respect to  $u^{\text{UCB}}$  by definition. This implies that  $s$  is feasible.

We see that the platform's revenue is equal to:

$$\begin{aligned} -\sum_{t=1}^T \sum_{a \in \mathcal{A}_t} \tau_a &= -\sum_{t=1}^T \sum_{a \in \mathcal{A}_t} \tau_a^* + \sum_{t=1}^T \sum_{a \in \mathcal{A}_t} \varepsilon + \sum_{t=1}^T \sum_{a \in \mathcal{A}_t} (\max(C_{a, \mu_X(a)}) - \min(C_{a, \mu_X(a)})) \\ &= \varepsilon \sum_{t=1}^T |\mathcal{A}_t| - \sum_{t=1}^T \sum_{a \in \mathcal{A}_t} (\max(C_{a, \mu_X(a)}) - \min(C_{a, \mu_X(a)})). \end{aligned}$$

Using the proof of Theorem 5.1, we see that

$$\sum_{t=1}^T \sum_{a \in \mathcal{A}_t} (\max(C_{a, \mu_X(a)}) - \min(C_{a, \mu_X(a)})) \leq O(|\mathcal{A}| \sqrt{nT} \log(|\mathcal{A}|T)),$$

as desired.  $\square$

## F PROOFS FOR SECTION 6.3

### F.1 Proof of Proposition 6.5

**PROOF OF PROPOSITION 6.5.** We first prove the first part of the statement, and then the second part of the statement.

*Proof of part (a).* We note that it follows immediately from Definition 6.4 that NTU Subset Instability is nonnegative. Let us now show that  $I(X; u, \mathcal{A})$  is zero if and only if  $(X, \tau)$  is stable. It is not difficult to see that the infimum of  $(\dagger)$  is attained at some  $s^*$ .

If  $I(X; u, \mathcal{A}) = 0$ , then we know that  $s_a^* = 0$  for all  $a \in \mathcal{A}$ . The constraints in the optimization problem imply that  $X$  has no blocking pairs and individual rationality is satisfied, as desired.

If  $X$  is stable, then we see that  $s = \vec{0}$  is a feasible solution to  $(\dagger)$ , which means that the optimum of  $(\dagger)$  is at most zero. This, coupled with the fact that  $I(X; u, \mathcal{A})$  is always nonnegative, means that  $I(X; u, \mathcal{A}) = 0$ , as desired.

---

**ALGORITHM 6:** COMPUTEMATCHNTU: Compute matching with transfers from confidence sets.
 

---

```

1: procedure COMPUTEMATCHNTU( $\mathcal{C}$ )
2:   for  $(i, j) \in \mathcal{I} \times \mathcal{J}$  do
3:      $u_i^{\text{UCB}}(j) \leftarrow \max(C_{i,j}); \quad u_j^{\text{UCB}}(i) \leftarrow \max(C_{j,i})$  ▷ UCB estimates of utilities.
4:   Run any version of the Gale-Shapley algorithm [23] on  $u^{\text{UCB}}$  to obtain a matching  $X^*$ .
5:   return  $X^*$ 

```

---

*Proof of part (b).* Consider two utility functions  $u$  and  $\tilde{u}$ . To show Lipchitz continuity, it suffices to show that for any matching  $X$ :

$$|I(X; u, \mathcal{A}) - I(X; \tilde{u}, \mathcal{A})| \leq 2 \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_\infty.$$

We show that:

$$I(X; \tilde{u}, \mathcal{A}) \leq I(X; u, \mathcal{A}) + 2 \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_\infty,$$

noting that the other direction follows from an analogous argument. Let  $s^*$  be an optimal solution to  $(\dagger)$  for the utilities  $u$ . Consider the solution  $s_a = s_a^* + 2\|u_a - \tilde{u}_a\|_\infty$ . We first verify that  $s$  is a feasible solution to  $(\dagger)$  for  $\tilde{u}$ . We see that:

$$\begin{aligned}
& \min(\tilde{u}_i(j) - \tilde{u}_i(\mu_X(i)) - s_i, \tilde{u}_j(i) - \tilde{u}_j(\mu_X(j)) - s_j) \\
&= \min(\tilde{u}_i(j) - \tilde{u}_i(\mu_X(i)) - s_i^* - 2\|u_i - \tilde{u}_i\|_\infty, \tilde{u}_j(i) - \tilde{u}_j(\mu_X(j)) - s_j^* - 2\|u_j - \tilde{u}_j\|_\infty) \\
&\leq \min(u_i(j) - u_i(\mu_X(i)) - s_i^*, u_j(i) - u_j(\mu_X(j)) - s_j^*) \\
&\leq 0,
\end{aligned}$$

as desired. Moreover, we see that

$$\tilde{u}_a(\mu_X(a)) + s_a = \tilde{u}_a(\mu_X(a)) + s_a^* + 2\|u_a - \tilde{u}_a\|_\infty \leq u_a(\mu_X(a)) + s_a^* \geq 0.$$

Thus, we have demonstrated that  $s$  is feasible. This means that:

$$I(X; \tilde{u}, \mathcal{A}) \leq \sum_{a \in \mathcal{A}} s_a = \sum_{a \in \mathcal{A}} s_a^* + 2 \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_\infty = [I(X; u, \mathcal{A}) + 2 \sum_{a \in \mathcal{A}} \|u_a - \tilde{u}_a\|_\infty],$$

as desired. □

## F.2 Proof of Theorem 6.6

We show that the algorithmic approach from Section 5 can be adapted to the setting of matching with non-transferable utilities.

Drawing intuition from Section 5, at each round, we compute a stable matching for utilities given by the upper confidence bounds. More precisely, suppose we have a collection  $\mathcal{C}$  of confidence sets  $C_{i,j}, C_{j,i} \subseteq \mathbb{R}$  such that  $u_i(j) \in C_{i,j}$  and  $u_j(i) \in C_{j,i}$  for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$ . Our algorithm uses  $\mathcal{C}$  to get an upper confidence bound for each agent's utility function and then computes a stable matching with transfers as if these upper confidence bounds were the true utilities (see COMPUTEMATCHNTU). This can be implemented efficiently if we use, e.g., the Gale-Shapley algorithm (either the customer-proposing algorithm or the provider-proposing algorithm will work).

The core property of COMPUTEMATCHNTU is that we can upper bound NTU Subset Instability by the sum of the sizes of the relevant confidence sets, assuming that the confidence sets contain the true utilities.

PROPOSITION F.1. Consider a collection confidence sets  $\mathcal{C}$  such that  $u_i(j) \in C_{i,j}$  and  $u_j(i) \in C_{j,i}$  for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$ . The instability of the output  $X^{\text{UCB}}$  of COMPUTEMATCH satisfies

$$I(X^{\text{UCB}}; u, \mathcal{A}) \leq \sum_{a \in \mathcal{A}^t} (\max(C_{a, \mu_{X^{\text{UCB}}}(a)}) - \min(C_{a, \mu_{X^{\text{UCB}}}(a)})). \quad (14)$$

PROOF. We construct subsidies for this setting to be:

$$s_a = \max(C_{a, \mu_X(a)}) - u_a(\mu_X(a)) \leq \max(C_{a, \mu_X(a)}) - \min(C_{a, \mu_X(a)}).$$

Step 1: Verifying feasibility. We first show that  $s$  is a feasible solution to  $(\dagger)$ .

$$\begin{aligned} & \min(u_i(j) - u_i(\mu_{X^{\text{UCB}}}(i)) - s_i, \tilde{u}_j(i) - u_j(\mu_{X^{\text{UCB}}}(j)) - s_j) \\ &= \min(u_i(j) - u_i^{\text{UCB}}(\mu_{X^{\text{UCB}}}(i)), \tilde{u}_j(i) - u_j^{\text{UCB}}(\mu_{X^{\text{UCB}}}(j))) \\ &\leq \min(u_i^{\text{UCB}}(j) - u_i^{\text{UCB}}(\mu_{X^{\text{UCB}}}(i)), u_j^{\text{UCB}}(i) - u_j^{\text{UCB}}(\mu_{X^{\text{UCB}}}(j))) \\ &\leq 0, \end{aligned}$$

where the last step uses the fact that  $\mu_{X^{\text{UCB}}}$  is stable with respect to  $u^{\text{UCB}}$  by definition. Moreover, we see that

$$u_a(\mu_{X^{\text{UCB}}}(a)) + s_a = u_a^{\text{UCB}}(\mu_{X^{\text{UCB}}}(a)) \geq 0,$$

where the last inequality uses that  $\mu_{X^{\text{UCB}}}$  is stable with respect to  $u^{\text{UCB}}$  by definition. This implies that  $s$  is feasible.

Step 2: Computing the objective. We next compute the objective of  $(\dagger)$  at  $s$  and use this to bound  $I(X^*; u, \mathcal{A})$ . A simple calculation shows that:

$$I(X^*; u, \mathcal{A}) \leq \sum_a s_a = \sum_{a \in \mathcal{A}} (\max(C_{a, \mu_{X^{\text{UCB}}}(a)}) - \min(C_{a, \mu_{X^{\text{UCB}}}(a)})),$$

as desired.  $\square$

**F.2.1 Explicit Algorithm and Regret Bounds.** Using the same intuition as Section 5, the regret bound of Proposition F.1 hints at an algorithm: Each round, select the matching with transfers returned by COMPUTEMATCHNTU and update confidence sets accordingly. To instantiate this approach, it remains to construct confidence intervals that contain the true utilities with high probability.

We showcase this algorithm in the simple setting of unstructured preferences. For this setting, we can construct our confidence intervals following the classical UCB approach. That is, for each utility value involving the pair  $(i, j)$ , we take a length  $O(\sqrt{\log(|\mathcal{A}|T)/n_{ij}})$  confidence interval centered around the empirical mean, where  $n_{ij}$  is the number of times the pair has been matched before. We describe this construction precisely in Algorithm 3 (MATCHNTUUCB).

To analyze MATCHNTUUCB, recall that Lemma 5.4 bounds the regret at each step by the lengths of the confidence intervals of each pair in the selected matching. Like in Section 5, this yields the following instance-independent regret bound:

$$\text{THEOREM F.2. } \text{MATCHNTUUCB incurs expected regret } \mathbb{E}(R_T) \leq O(|\mathcal{A}|^{3/2} \sqrt{T} \sqrt{\log(|\mathcal{A}|T)}).$$

PROOF. This proof proceeds very similarly to the proof of Theorem 5.1. We consider the event  $E$  that all of the confidence sets contain their respective true utilities at every timestep  $t \leq T$ . That is,  $u_i(j) \in C_{i,j}$  and  $u_j(i) \in C_{j,i}$  for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$  at all  $t$ .

**ALGORITHM 7:** MATCHNTUUCB: A bandit algorithm for matching with non-transferable utilities.

---

```

1: procedure MATCHNTUUCB( $T$ )
2:   for  $(i, j) \in \mathcal{I} \times \mathcal{J}$  do                                 $\triangleright$  Initialize confidence intervals and empirical mean.
3:      $C_{i,j} \leftarrow [-1, 1]; \quad C_{j,i} \leftarrow [-1, 1]; \quad \hat{u}_i(j) \leftarrow 0; \quad \hat{u}_j(i) \leftarrow 0$ 
4:   for  $1 \leq t \leq T$  do
5:      $X^t \leftarrow \text{COMPUTEMATCHNTU}(C)$ 
6:     for  $(i, j) \in X^t$  do                                 $\triangleright$  Set confidence intervals and update means.
7:       Update  $\hat{u}_i(j)$  and  $\hat{u}_j(i)$  from feedback; increment counter  $n_{ij}$ 
8:        $C_{i,j} \leftarrow [\hat{u}_i(j) - 8\sqrt{\log(|\mathcal{A}|T)/n_{ij}}, \hat{u}_i(j) + 8\sqrt{\log(|\mathcal{A}|T)/n_{ij}}] \cap [-1, 1]$ 
9:        $C_{j,i} \leftarrow [\hat{u}_j(i) - 8\sqrt{\log(|\mathcal{A}|T)/n_{ij}}, \hat{u}_j(i) + 8\sqrt{\log(|\mathcal{A}|T)/n_{ij}}] \cap [-1, 1]$ 

```

---

Case 1:  $E$  holds. By Lemma 5.4, we may bound

$$I(X^t; u, \mathcal{A}^t) \leq \sum_{a \in \mathcal{A}^t} (\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)})) = O\left(\sum_{(i,j) \in X^t} \sqrt{\frac{\log(|\mathcal{A}|T)}{n_{ij}^t}}\right),$$

where  $n_{ij}^t$  is the number of times that the pair  $(i, j)$  has been matched at the start of round  $t$ . Let  $w_{i,j}^t = \frac{1}{\sqrt{n_{ij}^t}}$  be the size of the confidence set (with the log factor scaled out) for  $(i, j)$  at the start of round  $t$ .

At each timestep  $t$ , let us consider the list consisting of  $w_{i_t, j_t}^t$  for all  $(i_t, j_t) \in X^t$ . Let us now consider the overall list consisting of the concatenation of all of these lists over all rounds. Let us order this list in decreasing order to obtain a list  $\tilde{w}_1, \dots, \tilde{w}_L$  where  $L = \sum_{t=1}^T |X^t| \leq nT$ . In this notation, we observe that:

$$\sum_{t=1}^T I(X^t; u, \mathcal{A}^t) \leq \sum_{t=1}^T \sum_{a \in \mathcal{A}^t} (\max(C_{a, \mu_{X^t}(a)}) - \min(C_{a, \mu_{X^t}(a)})) = \log(|\mathcal{A}|T) \sum_{l=1}^L \tilde{w}_l.$$

We claim that  $\tilde{w}_l \leq O(\min(1, \frac{1}{\sqrt{(l/|\mathcal{A}|^2)-1}}))$ . The number of rounds that a pair of agents can have their confidence set have size at least  $\tilde{w}_l$  is upper bounded by  $1 + \frac{1}{\tilde{w}_l^2}$ . Thus, the total number of times that any confidence set can have size at least  $\tilde{w}_l$  is upper bounded by  $(|\mathcal{A}|^2)(1 + \frac{1}{\tilde{w}_l^2})$ .

Putting this together, we see that:

$$\begin{aligned} \log(|\mathcal{A}|T) \sum_{l=1}^L \tilde{w}_l &\leq O\left(\sum_{l=1}^L \min\left(1, \frac{1}{\sqrt{(l/|\mathcal{A}|^2)-1}}\right)\right) \\ &\leq O\left(\log(|\mathcal{A}|T) \sum_{l=1}^{nT} \min\left(1, \frac{1}{\sqrt{(l/|\mathcal{A}|^2)-1}}\right)\right) \\ &\leq O(|\mathcal{A}|\sqrt{nT} \log(|\mathcal{A}|T)). \end{aligned}$$

Case 2:  $E$  does not hold. Since each  $n_{ij}(\hat{u}_i(j) - u_i(j))$  is mean-zero and 1-subgaussian, and we have  $O(|\mathcal{I}||\mathcal{J}|T)$  such random variables over  $T$  rounds, the probability that any of them exceeds

$$2\sqrt{\log(|\mathcal{I}||\mathcal{J}|T/\delta)} \leq 2\sqrt{\log(|\mathcal{A}|^2T/\delta)}$$



is at most  $\delta$  by a standard tail bound for the maximum of subgaussian random variables. It follows that  $E$  fails to hold with probability at most  $|\mathcal{A}|^{-2}T^{-2}$ . In the case that  $E$  fails to hold, our regret in any given round would be at most  $4|\mathcal{A}|$  by the Lipschitz property in Proposition 6.5. (Recall that our upper confidence bound for any utility is wrong by at most two due to clipping each confidence interval to lie in  $[-1, 1]$ .) Thus, the expected regret from this scenario is at most

$$|\mathcal{A}|^{-2}T^{-2} \cdot 4|\mathcal{A}|T \leq 4|\mathcal{A}|^{-1}T^{-1},$$

which is negligible compared to the regret bound from when  $E$  does occur.  $\square$

## REFERENCES

- [1] Jacob D. Abernethy, Elad Hazan, and Alexander Rakhlin. 2008. Competing in the dark: An efficient algorithm for bandit linear optimization. In *21st Annual Conference on Learning Theory (COLT'08)*. Omnipress, 263–274.
- [2] Max Alston. 2020. On the non-existence of stable matches with incomplete information. *Games Econ. Behav.* 120 (2020), 336–344.
- [3] Guy Aridor, Yishay Mansour, Aleksandrs Slivkins, and Zhiwei Steven Wu. 2020. Competing bandits: The perils of exploration under competition. *CoRR* abs/2007.10144 (2020).
- [4] Itai Ashlagi, Mark Braverman, Yash Kanoria, and Peng Shi. 2020. Clearing matching markets efficiently: Informative signals and match recommendations. *Manag. Sci.* 66, 5 (2020), 2163–2193.
- [5] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. 2002. Finite-time analysis of the multiarmed bandit problem. *Mach. Learn.* 47, 2–3 (2002), 235–256.
- [6] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. 2002. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.* 32, 1 (2002), 48–77.
- [7] Eduardo M. Azevedo and John William Hatfield. 2018. Existence of equilibrium in large matching markets with complementarities. Available at SSRN (2018). [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3268884](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3268884).
- [8] Ashwinkumar Badanidiyuru, Robert Kleinberg, and Aleksandrs Slivkins. 2018. Bandits with knapsacks. *J. ACM* 65, 3 (2018), 13:1–13:55.
- [9] Soumya Basu, Karthik Abinav Sankararaman, and Abishek Sankararaman. 2021. Beyond  $\log^2(T)$  regret for decentralized bandits in matching markets. *CoRR* abs/2103.07501 (2021).
- [10] Martin Bichler, Maximilian Fichtl, and Gregor Schwarz. 2021. Walrasian equilibria from an optimization perspective: A guide to the literature. *Naval Res. Logist.* 68, 4 (2021), 496–513.
- [11] Sushil Bikhchandani. 2017. Stability with one-sided incomplete information. *J. Econ. Theor.* 168 (2017), 372–399.
- [12] Sarah H. Cen and Devavrat Shah. 2021. Regret, stability, and fairness in matching markets with bandit learners. *CoRR* abs/2102.06246 (2021).
- [13] Nicolò Cesa-Bianchi and Gábor Lugosi. 2012. Combinatorial bandits. *J. Comput. Syst. Sci.* 78, 5 (2012), 1404–1422.
- [14] Chandra R. Chegireddy and Horst W. Hamacher. 1987. Algorithms for finding K-best perfect matchings. *Discret. Appl. Math.* 18, 2 (1987), 155–165.
- [15] Wei Chen, Yajun Wang, and Yang Yuan. 2013. Combinatorial multi-armed bandit: General framework and applications. In *30th International Conference on Machine Learning (JMLR Workshop and Conference Proceedings, Vol. 28)*. JMLR.org, 151–159.
- [16] Richard Combes, Mohammad Sadegh Talebi, Alexandre Proutière, and Marc Lelarge. 2015. Combinatorial bandits revisited. In *Annual Conference on Neural Information Processing Systems*. 2116–2124.
- [17] Sanmay Das and Emir Kamenica. 2005. Two-sided bandits and the dating market. In *19th International Joint Conference on Artificial Intelligence*. Professional Book Center, 947–952.
- [18] Gerard Debreu and Herbert Scarf. 1963. A limit theorem on the core of an economy. *Int. Econ. Rev.* 4, 3 (1963), 235–246.
- [19] Federico Echenique, Sangmok Lee, Matthew Shum, and M. Bumin Yenmez. 2013. The revealed preference theory of stable and extremal stable matchings. *Econometrica* 81, 1 (2013), 153–171.
- [20] Ehsan Emamjomeh-Zadeh, Yannai A. Gonczarowski, and David Kempe. 2020. The complexity of interactively learning a stable matching by trial and error. In *21st ACM Conference on Economics and Computation*. ACM.
- [21] Peter I. Frazier, David Kempe, Jon M. Kleinberg, and Robert Kleinberg. 2014. Incentivizing exploration. In *ACM Conference on Economics and Computation (EC'14)*. ACM, 5–22.
- [22] Yi Gai, Bhaskar Krishnamachari, and Rahul Jain. 2012. Combinatorial network optimization with unknown variables: Multi-armed bandits with linear rewards and individual observations. *IEEE/ACM Trans. Netw.* 20, 5 (2012), 1466–1478.
- [23] D. Gale and L. S. Shapley. 1962. College admissions and the stability of marriage. *Amer. Math. Month.* 69, 1 (1962), 9–15.

- [24] Yannai A. Gonczarowski, Noam Nisan, Rafail Ostrovsky, and Will Rosenbaum. 2019. A stable marriage requires communication. *Games Econ. Behav.* 118 (2019), 626–647.
- [25] Nicole Immorlica, Karthik Abinav Sankararaman, Robert E. Schapire, and Aleksandrs Slivkins. 2019. Adversarial bandits with knapsacks. In *60th IEEE Annual Symposium on Foundations of Computer Science (FOCS'19)*. IEEE Computer Society, 202–219.
- [26] Ramesh Johari, Vijay Kamble, and Yash Kanoria. 2021. Matching while learning. *Oper. Res.* 69, 2 (2021), 655–681.
- [27] Robert D. Kleinberg and Frank Thomson Leighton. 2003. The value of knowing a demand curve: Bounds on regret for online posted-price auctions. In *44th Symposium on Foundations of Computer Science (FOCS'03)*. IEEE Computer Society, 594–605.
- [28] H. W. Kuhn. 1955. The Hungarian method for the assignment problem. *Naval Res. Logist. Quart.* 2, 1–2 (1955), 83–97.
- [29] Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvári. 2015. Tight regret bounds for stochastic combinatorial semi-bandits. In *18th International Conference on Artificial Intelligence and Statistics (JMLR Workshop and Conference Proceedings, Vol. 38)*. JMLR.org.
- [30] Tor Lattimore and Csaba Szepesvári. 2020. *Bandit Algorithms*. Cambridge University Press.
- [31] Xiaocheng Li, Chunlin Sun, and Yinyu Ye. 2021. The symmetry between arms and knapsacks: A primal-dual approach for bandits with knapsacks. *CoRR* abs/2102.06385 (2021).
- [32] Lydia T. Liu, Horia Mania, and Michael I. Jordan. 2020. Competing bandits in matching markets. In *23rd International Conference on Artificial Intelligence and Statistics (Proceedings of Machine Learning Research, Vol. 108)*. PMLR, 1618–1628.
- [33] Lydia T. Liu, Feng Ruan, Horia Mania, and Michael I. Jordan. 2020. Bandit learning in decentralized matching markets. *CoRR* abs/2012.07348 (2020).
- [34] Qingmin Liu. 2020. Stability and Bayesian consistency in two-sided markets. *Amer. Econ. Rev.* 110, 8 (Aug. 2020).
- [35] Qingmin Liu, George J. Mailath, Andrew Postlewaite, and Larry Samuelson. 2014. Stable matching with incomplete information. *Econometrica* 82, 2 (2014), 541–587.
- [36] Yishay Mansour, Aleksandrs Slivkins, and Vasilis Syrgkanis. 2015. Bayesian incentive-compatible bandit exploration. In *16th ACM Conference on Economics and Computation (EC'15)*. ACM, 565–582.
- [37] J. E. Martínez-Legaz. 1996. Dual representation of cooperative games based on Fenchel-Moreau conjugation. *Optimization* 36, 4 (1996), 291–319.
- [38] Michael Rothschild. 1974. A two-armed bandit theory of market pricing. *J. Econ. Theor.* 9, 2 (1974), 185–202.
- [39] Daniel Russo and Benjamin Van Roy. 2013. Eluder dimension and the sample complexity of optimistic exploration. In *27th Annual Conference on Neural Information Processing Systems*. 2256–2264.
- [40] Abishek Sankararaman, Soumya Basu, and Karthik Abinav Sankararaman. 2021. Dominate or delete: Decentralized competing bandits in serial dictatorship. In *24th International Conference on Artificial Intelligence and Statistics (Proceedings of Machine Learning Research, Vol. 130)*. PMLR, 1252–1260.
- [41] L. S. Shapley and M. Shubik. 1966. Quasi-cores in a monetary economy with nonconvex preferences. *Econometrica* 34, 4 (1966), 805–827.
- [42] L. S. Shapley and M. Shubik. 1971. The assignment game I: The core. *Int. J. Game Theor.* 1, 1 (Dec. 1971), 111–130.
- [43] Peng Shi. 2020. Efficient matchmaking in assignment games with application to online platforms. In *21st ACM Conference on Economics and Computation*. ACM, 601–602.
- [44] Andrea Tirinzoni, Matteo Pirodda, Marcello Restelli, and Alessandro Lazaric. 2020. An asymptotically optimal primal-dual incremental algorithm for contextual linear bandits. In *Annual Conference on Neural Information Processing Systems*.

Received 31 December 2021; revised 24 August 2022; accepted 11 November 2022