

A Clarification of Notation

In this section, we provide a comprehensive clarification on the use of notation in this paper.

Throughout the paper, we use $\mathcal{O}(\cdot)$ to hide problem-independent constants and use $(\cdot)_{[a,b]}$ to denote the truncation into the range $[a, b]$. \mathcal{I} and \mathcal{J} denote the sets of all side-1 and side-2 agents respectively. Further, we use $I \in 2^{\mathcal{I}}$ and $J \in 2^{\mathcal{J}}$ to denote any set of participating agents, which are the subsets of \mathcal{I} and \mathcal{J} .

For any given stage $h \in [H]$, we use s_h and (C_h, I_h, J_h) interchangeably when describing any state $s_h \in \mathcal{S}$ where $\mathcal{S} = \mathcal{C} \times 2^{\mathcal{I}} \times 2^{\mathcal{J}}$. Analogously, we also use a_h and (e_h, X_h, τ_h) interchangeably for the action $a_h \in \mathcal{A}$ where $\mathcal{A} = \Upsilon \times \mathcal{X} \times \mathcal{T}$

We use $\pi = \{\pi_h\}_{h=1}^H$ to denote a policy, where each π_h is defined to be a mapping from \mathcal{S} to a distribution $\Delta_{\mathcal{A}}$ on \mathcal{A} . Therefore, for any $h \in [H]$ and $s_h \in \mathcal{S}$, $\pi_h(\cdot|s_h)$ denotes a probability distribution on \mathcal{A} . Note that because $\mathcal{A} = \Upsilon \times \mathcal{X} \times \mathcal{T}$, the policy π is a joint policy. We may slightly abuse the notation in the paper and refer to π as the policy restricted on Υ only, whenever it is clear from the context. In such case, we refer to $\pi_h(\cdot|s_h)$ as a distribution on Υ only.

We also present the following table of notations. The π in the superscript can be replaced by π_k or π_* , where the former refers to the policy in episode k , and the latter refers to the optimal policy.

Table 1: Notation

Notation	Meaning
$\mathcal{C}, \mathcal{I}, \mathcal{J}$	set of contexts, side-1 agents, side-2 agents
$\Upsilon, \mathcal{X}, \mathcal{T}$	set of planner's actions, all matchings and transfers over all possible subsets of $\mathcal{I} \times \mathcal{J}$
r_h, \bar{r}_h	reward, pseudo-reward functions
V_h^π, Q_h^π, V_h^*	value, Q functions under π , optimal value functions w.r.t. the transition functions $\{\mathbb{P}_h\}_{h=1}^H$ and reward functions $\{r_h\}_{h=1}^H$
$\bar{V}_h^\pi, \bar{Q}_h^\pi, \bar{V}_h^*$	pseudo-value, pseudo-Q functions under π , optimal pseudo-value functions w.r.t. the transition functions $\{\mathbb{P}_h\}_{h=1}^H$ and reward functions $\{\bar{r}_h\}_{h=1}^H$
\bar{V}_h^k, \bar{Q}_h^k	estimated value, Q functions for stage h in episode k in Algorithm 1
π_k	the policy followed by Algorithm 1 in episode k , where $\pi_k = \{\pi_{k,h}\}_{h=1}^H$
$\pi_{k,h}$	the policy followed by Algorithm 1 at stage h in episode k

B Supplementary Information on Matching and Stability

In this section, we review some basics on the matching problem. We first introduce the classic problem of (static) matching with transfers and the notion of stability. We then recap the primal-dual formulation that provides an efficient way to solve a stable matching (Shapley and Shubik, 1971). Finally, we give more details about Subset Instability and its properties.

B.1 Matching with Transferable Utilities

This section is a supplementary to Section 3.1 in the main text. We introduce the two-sided static matching with transferable utilities.

Denote the sets of participating agents by I and J for two sides respectively. A matching $X \subseteq I \times J$ is a set of pairs of agents, and $(i, j) \in X$ means $i \in I$ is matched to $j \in J$. Each agent is matched at most once. We denote by $X(i) = j$ and $X(j) = i$ for any matched pair $(i, j) \in X$, while for any unmatched agent $a \in I \cup J$, we write $X(a) = a$.

Matched agents receive utilities, denoted by $u : I \times J \rightarrow \mathbb{R}$ for agents in I and $v : I \times J \rightarrow \mathbb{R}$ for agents in J . Specifically, if $(i, j) \in X$, then agent i receives an utility $u(i, j)$ and agent j receives an

utility $v(i, j)$. Remaining unmatched agents receive zero utility. With a slight abuse of notation, we overwrite $u(i, i) = 0$ and $v(j, j) = 0$ for any $i \in I$ and $j \in J$.

In addition, there are utility transfers between (and only between) agents. We denote the transfer function by $\tau : I \times J \rightarrow \mathbb{R}$ such that for any agent $a \in I \cup J$, $\tau(a)$ is the transfer received by agent a . Since the transfers are within agents, we have

$$\sum_{a \in I \cup J} \tau(a) = 0.$$

We denote the market outcome by (X, τ) , under which the net utility received by an agent $i \in I$ is $u(i, j) + \tau(i)$ if $(i, j) \in X$, and similarly for agents in J .

The notion of *stable matching* is as follows.

Definition B.1 (Stable matching). A matching-transfer pair (X, τ) on I, J is stable if:

1. The net utility of of any agent is non-negative, i.e.

$$\begin{aligned} u(i, X(i)) + \tau(i) &\geq 0, \\ v(X(j), j) + \tau(j) &\geq 0, \end{aligned}$$

for all $i \in I$ and $j \in J$.

2. There are no blocking pairs, i.e.

$$[u(i, X(i)) + \tau(i)] + [v(X(j), j) + \tau(j)] \geq u(i, j) + v(i, j),$$

for all pairs $(i, j) \in I \times J$.

Stable matching implies that no matched agents would rather be unmatched and no pair of agents can find a transfer between themselves so that both would rather match with each other than follow (X, τ) . The following proposition provides a fundamental and important max-weight interpretation for stable matchings.

Proposition B.2 (Shapley and Shubik 1971). For the matching with transfer problem, if (X, τ) is a stable matching under Definition B.1 then X must be the max-weight matching, i.e.,

$$X = \arg \max_{X'} \sum_{i \in I, j \in J} u(i, X'(i)) + v(X'(j), j)$$

where the maximum is over all matchings on $I \times J$.

Therefore, by Proposition B.2 to maximize the total social welfare (i.e. sum of utilities), it suffices to find a stable matching. But how? This is answered in the next subsection.

B.2 The Linear Program and Dual Program

In this subsection, we explain how to find a stable matching (X, τ) given input I, J, u, v , which gives rise to the algorithm OM (i.e. Algorithm 3) in the main text.

Shapley and Shubik (1971) showed that, assuming the utility functions are known, the stable (X, τ) can be found by solving the following linear program and its dual program (recapped from Section 4.1 in the main text):

$$\begin{aligned} \mathcal{LP}(I, J, u, v) : \quad & \max_{w \in \mathbb{R}^{|I| \times |J|}} \sum_{(i,j) \in I \times J} w_{i,j} [u(i, j) + v(i, j)] \\ & \text{s.t.} \quad \sum_{j \in J_h} w_{i,j} \leq 1, \forall i \in I, \\ & \quad \sum_{i \in I_h} w_{i,j} \leq 1, \forall j \in J, \\ & \quad w_{i,j} \geq 0, \forall (i, j) \in I \times J, \end{aligned} \tag{14}$$

and its dual program:

$$\begin{aligned} \mathcal{DP}(I, J, u, v) : \quad & \min_{p: I \cup J \rightarrow \mathbb{R}^+} \sum_{a \in I \cup J} p(a) \\ \text{s.t.} \quad & p(i) + p(j) \geq u(i, j) + v(i, j), \forall (i, j) \in I \times J. \end{aligned} \quad (15)$$

Denote the solution pair to the primal-dual problems by (w, p) . [Shapley and Shubik \(1971\)](#) proved that (w, p) leads to a max-weight stable matching-transfer pair (X, τ) . Specifically, it is proved that the vector w must have integer entries, i.e., $w_{i,j} = 0$ or 1, which naturally induces a matching X such that $(i, j) \in X$ if and only if $w_{i,j} = 1$. Correspondingly, the transfers are $\tau(i) = p(i) - u(i, X(i))$ for $i \in I$, and similarly for $j \in J$.

The above procedure constitutes the subroutine oracle OM as displayed in [Algorithm 3](#). It takes as input the sets of participating agents and estimated utility functions, then outputs the stable matching (X, τ) by solving the primal-dual linear program described above. Note that the matching is only stable with respect to the estimated utility functions.

B.3 Details about Subset Instability

Next, we review the notion of Subset Instability and its properties. We refer the interested reader to [\(Jagadeesan et al. 2021\)](#) for the full details.

Given a matching-pair (X, τ) , define its utility difference as

$$\left[\max_{X'} \sum_{i \in I, j \in J} u(i, X'(i)) + v(X'(j), j) \right] - \left[\sum_{i \in I, j \in J} u(i, X(i)) + v(X(j), j) \right]. \quad (16)$$

Recall the definition of Subset Instability:

Definition 4.1 (Subset Instability, [Jagadeesan et al. 2021](#)). *Given any agent sets I, J and utility functions $u, v : I \times J \rightarrow \mathbb{R}$, the Subset Instability $\text{SI}(X, \tau; I, J, u, v)$ of the matching and transfer (X, τ) is defined as*

$$\begin{aligned} \max_{I' \times J' \subseteq I \times J} & \left[\left(\max_{X'} \sum_{i \in I'} u(i, X'(i)) + \sum_{j \in J'} v(X'(j), j) \right) \right. \\ & \left. - \left(\sum_{i \in I'} (u(i, X(j)) + \tau(i)) \right) - \left(\sum_{j \in J'} (v(X(j), j) + \tau(j)) \right) \right] \end{aligned}$$

where $X(\cdot)$ and $X'(\cdot)$ denotes the matched agent in matching X and X' respectively.

By definition, Subset Instability indicates whether there exists any subset $I' \times J'$ of agents who can achieve a higher total utility by taking some alternative matching X' among themselves other than the current matching-transfer pair (X, τ) . Thus, Subset Instability is an upper bound of the total utility difference, and quantifies the distance from a proposed matching to the optimal stable matching, as summarized in the following proposition.

Proposition B.3 (Proposition 4.4 in [Jagadeesan et al. 2021](#)). *The following holds for Subset Instability*

1. *Subset Instability is always nonnegative and is zero if and only if (X, τ) is stable matching.*
2. *Subset Instability is Lipschitz continuous with respect to the ℓ_∞ norm of the utility functions.*

$$\begin{aligned} & |\text{SI}(X, \tau; I, J, u, v) - \text{SI}(X, \tau; I, J, \tilde{u}, \tilde{v})| \\ & \leq 2 \left(\sum_{i \in I} \|u(i, \cdot) - \tilde{u}(i, \cdot)\|_\infty + \sum_{j \in J} \|v(\cdot, j) - \tilde{v}(\cdot, j)\|_\infty \right). \end{aligned}$$

3. *Subset Instability is always at least the utility difference [\(16\)](#).*

In our problem, this allows us to bound the total regret of the agents by the sum of Subset Instability of matchings across all episodes, as reflected in [Proposition 4.2](#).

C Proof Sketch

C.1 Optimistic Utility and Reward Estimates

The key step in our analysis is to show that the estimated pseudo-reward function \bar{r}_h^k satisfies optimism, i.e., $\bar{r}_h^k \geq \bar{r}_h$, and we need to ensure that \bar{r}_h^k is not too far away from \bar{r}_h . The lemma below justifies the optimism of utility estimates.

Lemma C.1 (UCB for Utility Estimates; proof in [E.1](#)). *For any $0 < \delta < 1$, set β_u as $\beta_u = \sqrt{d^2 \log[2(1 + d^2 K \max_h \min\{|I_h|, |J_h|\})/(\lambda\delta)]} + \sqrt{\lambda d}$. Then with probability at least $1 - \delta$, u_h^k and v_h^k in Algorithm [1](#) satisfy $u_h^k \geq u_h$ and $v_h^k \geq v_h$. Furthermore, $|u_h^k(\cdot) - u_h(\cdot)|$ and $|v_h^k(\cdot) - v_h(\cdot)|$ are bounded by $2\beta_u \|\Phi(\cdot)\|_{(\Sigma_h^k)^{-1}}$.*

Next, we explain why the optimism of utility estimates implies that of reward estimates. By Lemma [C.1](#), we can write $u_h^k = u_h + b_{u,h}$ and $v_h^k = v_h + b_{v,h}$ where $b_{u,h}$ and $b_{v,h}$ are bonus functions satisfying

$$b_{u,h}(C, e, i, j), b_{v,h}(C, e, i, j) \in [0, 2\beta_u \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}}].$$

Lemma C.2 (Planner's Optimism; proof in [E.2](#)). *Under the event of Lemma [C.1](#) it holds that for any $(C, e) \in \mathcal{C} \times \Upsilon$,*

$$0 \leq \bar{r}_h^k(C, e, I_h, J_h) - \bar{r}_h(C, e, I_h, J_h) \leq \sum_{(i,j) \in X_h^k} (b_{u,h}(C, e, i, j) + b_{v,h}(C, e, i, j)).$$

In the sequel, we denote by π_k the policy whose market making part is greedy w.r.t. \bar{Q}_h^k and whose matching part chooses the max-weight stable matching given u_h^k and v_h^k .

C.2 Proof Sketch of Theorem [5.4](#)

By definition [\(9\)](#), the agents' regret can be interpreted as the expected sum of total SI across all time steps, where the expectation is over the trajectory induced by π_k for $k \in [K]$.

To bound the regret, we relate the expected SI in [\(9\)](#) with the realized SI via a martingale difference sequence. Specifically, writing $\text{SI}_h = \text{SI}(s_h, a_h, u_h, v_h)$ and $\text{SI}_h^k = \text{SI}(s_h^k, a_h^k, u_h^k, v_h^k)$, we define the sum of differences as

$$\sum_{k=1}^K \left\{ \mathbb{E}_{\pi_k} \left[\sum_{h=1}^H \text{SI}_h \right] - \sum_{h=1}^H \text{SI}_h^k \right\}. \quad (17)$$

We bound the difference [\(17\)](#) and the sum of realized SI $\sum_{k=1}^K \sum_{h=1}^H \text{SI}_h^k$ separately, where the former is a sum of martingale difference sequences that concentrates and the latter can be bounded using the following lemma.

Lemma C.3 (Lemma 5.4 in [Jagadeesan et al. 2021](#); proof in [E.3](#)). *Under the event of Lemma [C.1](#) we have*

$$\text{SI}_h^k \leq \sum_{(i,j) \in X_h^k} (b_{u,h}(C_h^k, e_h^k, i, j) + b_{v,h}(C_h^k, e_h^k, i, j)).$$

Remark C.4. Note that each implemented matching induces several utility observations at a time, so bounding the bonus sum for utilities has a resemblance to lazy policy updates in the online learning literature ([Abbasi-Yadkori et al. 2011](#)). It particularly is similar to techniques used in the low switching cost problem in RL ([Bai et al. 2019](#); [Wang et al. 2021](#); [Gao et al. 2021](#)). This will be clear in the proof of Theorem [5.4](#) in Appendix [D.3](#).

C.3 Proof Sketch of Theorem [5.5](#)

Define the following functions δ_h^k and terms $\zeta_{k,h}^1, \zeta_{k,h}^2$:

$$\begin{aligned} \delta_h^k(C, e, I, J) &:= [\bar{r}_h + \mathbb{P}_h \bar{V}_{h+1}^k - \bar{Q}_h^k](C, e, I, J), \\ \zeta_{k,h}^1 &:= (\bar{V}_h^k - \bar{V}_h^{\pi_k})(C_h^k, I_h, J_h) - (\bar{Q}_h^k - \bar{Q}_h^{\pi_k})(C_h^k, e_h^k, I_h, J_h), \\ \zeta_{k,h}^2 &:= \mathbb{P}_h(\bar{V}_{h+1}^k - \bar{V}_{h+1}^{\pi_k})(C_h^k, e_h^k, I_h, J_h) - (\bar{V}_{h+1}^k - \bar{V}_{h+1}^{\pi_k})(C_{h+1}^k, I_{h+1}, J_{h+1}). \end{aligned} \quad (18)$$

To simplify the notation, in the following, we omit I_h, J_h from the arguments of the functions since we are conditioning on $\{I_h, J_h\}_{h=1}^H$ being fixed.

Lemma C.5 (Regret Decomposition of Planner; proof in [E.4](#)). *The planner's regret defined by [\(8\)](#) satisfies*

$$R^P(K) = \underbrace{\sum_{k=1}^K \sum_{h=1}^H (\zeta_{k,h}^1 + \zeta_{k,h}^2)}_{E_1} + \underbrace{\sum_{k=1}^K \sum_{h=1}^H [\mathbb{E}_{\pi^*} [\delta_h^k(C_h, e_h) | C_1^k] - \delta_h^k(C_h^k, e_h^k)]}_{E_2} \quad (19)$$

$$+ \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{\pi^*} \left[\langle \bar{Q}_h^k(C_h, \cdot), \pi_h^*(\cdot | C_h) - \pi_{k,h}(\cdot | C_h) \rangle_{\Upsilon} \right]}_{E_3} \Big| C_1^k,$$

where the expectation $\mathbb{E}_{\pi^*}[\cdot | C_1^k]$ is with respect to the trajectory $\{C_h, e_h\}_{h=1}^H$ induced by the policy π^* conditioning on $C_1 = C_1^k$ and $\langle \cdot, \cdot \rangle_{\Upsilon}$ means sum over all $e \in \Upsilon$.

In decomposition [\(19\)](#), term E_1 is controlled using a standard martingale concentration. Next, to bound E_2 , we show that $\delta_h^k \leq 0$ with high probability, which implies that $E_2 \leq \sum_{k=1}^K \sum_{h=1}^H |\delta_h^k(C_h^k, e_h^k)|$. Bounding each $|\delta_h^k(C_h^k, e_h^k)|$ by the corresponding optimistic bonus $\|\psi(C_h^k, e_h^k)\|_{(\Lambda_h^{k+1})^{-1}}$, we then apply Elliptical Potential lemma to get

$$E_2 \leq 2\beta_V \sum_{h=1}^H \sum_{k=1}^K \sqrt{2} \|\psi(C_h^k, e_h^k)\|_{(\Lambda_h^{k+1})^{-1}} \leq 2\sqrt{2}\beta_V H \sqrt{Kd \log((K+d)/d)}.$$

Finally, for E_3 , note that by Algorithm [1](#), the market-making part of the policy $\pi_{k,h}$ is the greedy policy with respect to \bar{Q}_h^k . Based on this observation, it follows that $\sum_{e \in \Upsilon} \bar{Q}_h^k(C_h, e) (\pi_h^*(e | C_h) - \pi_{k,h}(e | C_h)) = \sum_{e \in \Upsilon} \bar{Q}_h^k(C_h, e) \pi_h^*(e | C_h) - \max_{e \in \Upsilon} \bar{Q}_h^k(C_h, e) \leq 0$, so $E_3 \leq 0$. Combining yields the bound in Theorem [5.5](#). Full details are presented in Appendix [C.4](#). Note that our notion of meta algorithm for Algorithm [1](#) is different from that of meta learning ([Finn et al., 2017](#); [Xu et al., 2021](#)).

C.4 Proof of Theorem [5.5](#)

To prove Theorem [5.5](#), we need the following two lemmas which is helpful for bounding E_1 and E_2 .

Lemma C.6 (Proof in Section [E.5](#)). *Under the setting of Theorem [5.5](#) with probability at least $1 - \delta$, for all $(h, k) \in [H] \times [K]$ and $(C, e) \in \mathcal{C} \times \Upsilon$, it holds that*

$$-2\beta_V \cdot \|\psi(C, e)\|_{(\Lambda_h^k)^{-1}} \leq \delta_h^k(C, e) \leq 0.$$

Lemma C.7 (Proof in Section [E.6](#)). *For any $\delta > 0$, with probability at least $1 - \delta$, it holds that*

$$\sum_{k=1}^K \sum_{h=1}^H (\zeta_{k,h}^1 + \zeta_{k,h}^2) \leq 3 \left(\sum_{h=1}^H W_h \right) \sqrt{K \log 2/\delta}.$$

We are now ready to prove Theorem [5.5](#).

Proof of Theorem [5.5](#) By Lemma [C.5](#), we bound the three terms separately.

Bound on E_3 . According to Algorithm [1](#), the planner's policy $\pi_{k,h}$ is the greedy policy with respect to \bar{Q}_h^k . It follows that

$$\begin{aligned} & \langle \bar{Q}_h^k(C_h, I_h, J_h, \cdot), \pi_h^*(\cdot | C_h, I_h, J_h) - \pi_{k,h}(\cdot | C_h, I_h, J_h) \rangle_{\Upsilon} \\ &= \langle \bar{Q}_h^k(C_h, I_h, J_h, \cdot), \pi_h^*(\cdot | C_h, I_h, J_h) \rangle_{\Upsilon} - \max_{e \in \Upsilon} \bar{Q}_h^k(C_h, I_h, J_h, e) \leq 0. \end{aligned}$$

Therefore, we have $E_3 \leq 0$.

Bound on E_1 . We apply Lemma C.6. The E_1 term can be bounded as

$$E_1 \leq \sum_{k=1}^K \sum_{h=1}^H -\delta_h^k(C_h^k, e_h^k) \leq 2\beta_V \cdot \sum_{k=1}^K \sum_{h=1}^H \|\psi(C_h^k, e_h^k)\|_{(\Lambda_h^k)^{-1}}.$$

Note that $\Lambda_h^k = \lambda \mathbf{I} + \sum_{t=1}^{k-1} \psi(C_h^t, e_h^t) \psi(C_h^t, e_h^t)^\top$ by definition, where $\|\psi\|_2 \leq 1$ by Assumption 5.2 and $\lambda = 1$ by our choice. Thus we have $\Lambda_h^{k+1} = \Lambda_h^k + \psi(C_h^k, e_h^k) \psi(C_h^k, e_h^k)^\top \preceq 2\Lambda_h^k$, or equivalently, $(\Lambda_h^k)^{-1} \preceq 2(\Lambda_h^{k+1})^{-1}$, for all k . Applying this to the above inequality, we get the final bound for E_1 :

$$E_1 \leq 2\beta_V \cdot \sum_{h=1}^H \sum_{k=1}^K \sqrt{2} \|\psi(C_h^k, e_h^k)\|_{(\Lambda_h^{k+1})^{-1}} \leq 2\sqrt{2}\beta_V H \cdot \sqrt{Kd \cdot \log\left(\frac{K+d}{d}\right)}, \quad (20)$$

where the second step is by the Elliptical Potential Lemma (Lemma G.2).

Bound on E_2 . By Lemma C.7 we have

$$E_2 \leq 3 \left(\sum_{h=1}^H W_h \right) \cdot \sqrt{K \cdot \log \frac{2}{\delta}}. \quad (21)$$

Combining Lemma C.5, (20), (21) and $E_3 \leq 0$, we get that, with probability at least $1 - 3\delta$,

$$R^P(K) \leq 6\eta d^{5/2} H \left(\sum_{h=1}^H W_h \right) \cdot \sqrt{K} \cdot \log \left(\frac{dKH \min\{|\mathcal{I}|, |\mathcal{J}|\}}{\delta} \right),$$

where the $1 - 3\delta$ probability is from the union bound on the events of Lemma C.1, Lemma C.6 and Lemma C.7. Since $dKH \min\{|\mathcal{I}|, |\mathcal{J}|\}/\delta > 3$, replacing δ with $\delta/3$ and absorbing the constant into the big-O notation, we finish the proof. \square

D Proof of the Main Theory

D.1 Proof of Proposition 4.2

Proof of Proposition 4.2 By (3), the total regret can be written as

$$R(K) = \sum_{k=1}^K [V_1^*(s) - V_1^{\pi_k}(s)] = \sum_{k=1}^K [V_1^*(s) - \bar{V}_1^{\pi_k}(s)] + \sum_{k=1}^K [\bar{V}_1^{\pi_k}(s) - V_1^{\pi_k}(s)],$$

where $\bar{V}_1^{\pi_k}$ is the pseudo-value function defined by (7) corresponding to the pseudo-reward \bar{r}_h (defined by (6)) and induced by the policy π_k . Note that here we only care about the Υ part of π_k since matching-transfer has been maximized out by the definition of \bar{r}_h .

Now for any policy $\pi \in \Pi$ where $\pi = \{\pi_h\}_{h \in [H]}$, there exists a counterpart $\pi' = \{\pi'_h\}_{h \in [H]}$, such that $\pi_h(C|s) = \pi'_h(C|s)$ for all $C \in \Upsilon$ and $s \in \mathcal{S}$, whereas for the matching part π' always chooses the stable matching w.r.t. the true and unknown utility functions $u_h(C_h, e_h, \cdot, \cdot)$ and $v_h(C_h, e_h, \cdot, \cdot)$. Since π' can be viewed as a function of π , we write $\pi' = \pi'(\pi)$. Since the matching does not affect the context transition by assumption, and the stable matching maximizes total utility by Proposition B.2 we then have that

$$\begin{aligned} V_1^* &= \max_{\pi \in \Pi} V_1^\pi \\ &= \max_{\pi \in \Pi} \mathbb{E}_\pi \left[\sum_{h=1}^H r_h(s_h, a_h) \mid s_1 = s; a_h \sim \pi_h(\cdot | s_h), s_{h+1} \sim \mathbb{P}_h(\cdot | s_h, a_h), \forall h \in [H] \right] \\ &= \max_{\pi'(\pi): \pi \in \Pi} \mathbb{E}_{\pi'} \left[\sum_{h=1}^H r_h(s_h, a_h) \mid s_1 = s; a_h \sim \pi'_h(\cdot | s_h), s_{h+1} \sim \mathbb{P}_h(\cdot | s_h, a_h), \forall h \in [H] \right] \\ &= \max_{\pi'(\pi): \pi \in \Pi} \mathbb{E}_{\pi'} \left[\sum_{h=1}^H \bar{r}_h(s_h, e_h) \mid s_1 = s; e_h \sim \pi'_h(\cdot | s_h), s_{h+1} \sim \mathbb{P}_h(\cdot | s_h, e_h), \forall h \in [H] \right] = \bar{V}_1^*, \end{aligned}$$

where the fourth step is by definition of \bar{r}_h . It follows that

$$R(K) = \sum_{k=1}^K \left[\bar{V}_1^*(s) - \bar{V}_1^{\pi_k}(s) \right] + \sum_{k=1}^K \left[\bar{V}_1^{\pi_k}(s) - V_1^{\pi_k}(s) \right] = R^P(K) + \sum_{k=1}^K \left[\bar{V}_1^{\pi_k}(s) - V_1^{\pi_k}(s) \right]. \quad (22)$$

The second term in the R.H.S. can be written as

$$\begin{aligned} & \sum_{k=1}^K \left[\bar{V}_1^{\pi_k}(s) - V_1^{\pi_k}(s) \right] \\ &= \mathbb{E}_{\pi_k} \left[\sum_{h=1}^H \bar{r}_h(s_h, e_h) - r_h(s_h, a_h) \mid s_1 = s; a_h \sim \pi_{k,h}(\cdot | s_h), s_{h+1} \sim \mathbb{P}_h(\cdot | s_h, a_h), \forall h \in [H] \right]. \end{aligned}$$

Note that $\bar{r}_h(s_h, e_h) - r_h(s_h, a_h)$ is exactly the utility difference defined by (16). Since Subset Instability is at least the utility difference by Proposition B.3, it follows that

$$\begin{aligned} & \sum_{k=1}^K \left[\bar{V}_1^{\pi_k}(s) - V_1^{\pi_k}(s) \right] \\ & \leq \sum_{k=1}^K \mathbb{E}_{\pi_k} \left[\sum_{h=1}^H \text{SI}(s_h, a_h, u_h, v_h) \mid s_1 = s; a_h \sim \pi_{k,h}(\cdot | s_h), s_{h+1} \sim \mathbb{P}_h(\cdot | s_h, a_h), \forall h \in [H] \right] \\ & = R^M(K), \end{aligned}$$

where the last step is by the definition of $R^M(K)$ in (9). Plugging into (22), we get

$$R(K) \leq R^P(K) + R^M(K).$$

This completes the proof. \square

D.2 Bounds for Regression Estimators

We first bound the norm of the regression estimators in the algorithm.

Lemma D.1. *The regression estimators \mathbf{w}_h^k in Algorithm 1 satisfy*

$$\|\mathbf{w}_h^k\|_2 \leq \left(\sum_{l=h}^H W_l \right) \cdot \sqrt{\frac{dk}{\lambda}}, \quad \|\boldsymbol{\theta}_h^k\|_2 \leq \sqrt{\frac{d^2 k \cdot \min\{|\mathcal{I}|, |\mathcal{J}|\}}{\lambda}}, \quad \|\boldsymbol{\gamma}_h^k\|_2 \leq \sqrt{\frac{d^2 k \cdot \min\{|\mathcal{I}|, |\mathcal{J}|\}}{\lambda}}.$$

Proof of Lemma D.1 Consider \mathbf{w}_h^k for arbitrary h, k . For any vector $\mathbf{v} \in \mathbb{R}^d$,

$$\begin{aligned} |\mathbf{v}^\top \mathbf{w}_h^k| &= \left| (\boldsymbol{\Lambda}_h^k)^{-1} \sum_{t=1}^{k-1} \boldsymbol{\psi}(C_h^t, e_h^t) \bar{V}_{h+1}^k(C_{h+1}^t) \right| \\ &\leq \sum_{t=1}^{k-1} \left| \mathbf{v}^\top (\boldsymbol{\Lambda}_h^k)^{-1} \boldsymbol{\psi}(C_h^t, e_h^t) \right| \cdot \sum_{l=h}^H W_l \\ &\leq \left(\sum_{l=h}^H W_l \right) \cdot \sqrt{\left[\sum_{t=1}^{k-1} \mathbf{v}^\top (\boldsymbol{\Lambda}_h^k)^{-1} \mathbf{v} \right] \cdot \left[\sum_{t=1}^{k-1} \boldsymbol{\psi}(C_h^t, e_h^t)^\top (\boldsymbol{\Lambda}_h^k)^{-1} \boldsymbol{\psi}(C_h^t, e_h^t) \right]} \\ &\leq \left(\sum_{l=h}^H W_l \right) \cdot \sqrt{dk/\lambda} \cdot \|\mathbf{v}\|_2, \end{aligned}$$

where the first inequality holds because of the truncation of \bar{Q}_h^k and $\bar{V}_h^k(C, I, J) = \max_e \bar{Q}_h^k(C, e, I, J)$, the second inequality is by the Cauchy-Schwarz inequality, and the last inequality is by Lemma G.1. Since above holds for any \mathbf{v} , we conclude that

$$\|\mathbf{w}_h^k\|_2 = \max_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} |\mathbf{v}^\top \mathbf{w}_h^k| \leq \left(\sum_{l=h}^H W_l \right) \cdot \sqrt{dk/\lambda}.$$

For θ_h^k , let $\mathbf{v} \in \mathbb{R}^{d \times d}$. Then by the same analysis we have

$$\begin{aligned}
|\mathbf{v}^\top \theta_h^k| &= \left| \mathbf{v}^\top (\Sigma_h^k)^{-1} \sum_{t=1}^{k-1} \sum_{(i,j) \in X_h^t} \Phi(C_h^t, e_h^t, i, j) u_h^t(i, j) \right| \\
&\leq \sum_{t=1}^{k-1} \sum_{(i,j) \in X_h^t} \left| \mathbf{v}^\top (\Sigma_h^k)^{-1} \Phi(C_h^t, e_h^t, i, j) \right| \cdot 1 \\
&\leq \sqrt{\left[\sum_{t=1}^{k-1} \sum_{(i,j) \in X_h^t} \mathbf{v}^\top (\Sigma_h^k)^{-1} \mathbf{v} \right]} \cdot \sqrt{\left[\sum_{t=1}^{k-1} \sum_{(i,j) \in X_h^t} \Phi(C_h^t, e_h^t, i, j)^\top (\Sigma_h^k)^{-1} \Phi(C_h^t, e_h^t, i, j) \right]} \\
&\leq \|\mathbf{v}\|_2 \cdot \sqrt{k \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot d^2/\lambda},
\end{aligned}$$

which implies

$$\|\theta_h^k\|_2 \leq \sqrt{\frac{d^2 k \cdot \min\{|\mathcal{I}|, |\mathcal{J}|\}}{\lambda}}.$$

The same holds for γ_h^k . □

D.3 Proof of Theorem 5.4

We present the complete proof of Theorem 5.4

Proof of Theorem 5.4 Recall from Proposition B.3 that Subset Instability is at least the utility difference. We thus have

$$\begin{aligned}
R^M(K) &= \sum_{k=1}^K \left[\max_{\pi \in \Pi} V_1^\pi(s_1) - V_1^{\pi^k}(s_1) \right] \\
&\leq \sum_{k=1}^K \mathbb{E}_{\pi_k} \left[\sum_{h=1}^H \text{SI}(s_h, a_h, u_h, v_h) \right] \\
&= \sum_{k=1}^K \left\{ \underbrace{\mathbb{E}_{\pi_k} \left[\sum_{h=1}^H \text{SI}(s_h, a_h, u_h, v_h) \right]}_{R_k^M} - \underbrace{\sum_{h=1}^H \text{SI}(s_h^k, a_h^k, u_h, v_h)}_{\tilde{R}_k^M} \right\} + \sum_{k=1}^K \sum_{h=1}^H \text{SI}(s_h^k, a_h^k, u_h, v_h).
\end{aligned}$$

For each $k \in [K]$, let \mathcal{F}_{k-1} denote all the history until the beginning of Episode k . Then $\pi_k \sim \mathcal{F}_{k-1}$ and $R_k^M \sim \mathcal{F}_{k-1}$. Furthermore, we have $R_k^M = \mathbb{E}[\tilde{R}_k^M \mid \mathcal{F}_{k-1}]$, i.e., R_k^M is the conditional expectation of the realized quantity \tilde{R}_k^M . We can view (17) as the sum of a martingale difference sequence. From Definition 4.1, it holds almost surely that

$$|\mathbb{E}_{\pi_k} [\text{SI}(s_h, a_h, u_h, v_h)] - \text{SI}(s_h^k, a_h^k, u_h, v_h)| \leq W_h.$$

By Lemma G.3, we have that, for any $0 < \delta < 1$, with probability at least $1 - \delta$,

$$\sum_{l=1}^K (R_k^M - \tilde{R}_k^M) \leq \left(\sum_{h=1}^H W_h \right) \sqrt{2K \log \left(\frac{2}{\delta} \right)}. \quad (23)$$

To bound the second term $\sum_{k=1}^K \tilde{R}_k^M$, note that by Lemma C.1 and Lemma C.3, we have

$$\begin{aligned}
\sum_{k=1}^K \tilde{R}_k^M &\leq \sum_{k=1}^K \sum_{h=1}^H \sum_{(i,j) \in X_h^k} 4\beta_u \|\Phi(C_h^k, e_h^k, i, j)\|_{(\Sigma_h^k)^{-1}} \\
&= 4\beta_u \sum_{h=1}^H \left(\sum_{k=1}^K \sum_{(i,j) \in X_h^k} \|\Phi(C_h^k, e_h^k, i, j)\|_{(\Sigma_h^k)^{-1}} \right). \quad (24)
\end{aligned}$$

Note that by definition, we have

$$\Sigma_h^{k+1} = \Sigma_h^k + \sum_{(i,j) \in X_h^k} \Phi(C_h^k, e_h^k, i, j) \Phi(C_h^k, e_h^k, i, j)^\top.$$

By picking $\lambda = 1$ and the assumption that $\|\Phi\| \leq 1$, we have

$$\sum_{(i,j) \in X_h^k} \Phi(C_h^k, e_h^k, i, j) \Phi(C_h^k, e_h^k, i, j)^\top \preceq \min(|I_h|, |J_h|) \cdot \mathbf{I}_{d^2} \preceq \min(|I_h|, |J_h|) \cdot \Sigma_h^k.$$

It follows that $\Sigma_h^{k+1} \preceq (1 + \min(|I_h|, |J_h|)) \cdot \Sigma_h^k$, and thus $(\Sigma_h^k)^{-1} \preceq (1 + \min(|I_h|, |J_h|)) \cdot (\Sigma_h^{k+1})^{-1}$. Combining with (24), we get that

$$\begin{aligned} \sum_{k=1}^K \tilde{R}_k^M &\leq 4\beta_u \sum_{h=1}^H \left(\sqrt{1 + \min(|I_h|, |J_h|)} \cdot \sum_{k=1}^K \sum_{(i,j) \in X_h^k} \|\Phi(C_h^k, e_h^k, i, j)\|_{(\Sigma_h^{k+1})^{-1}} \right) \\ &\leq 4\sqrt{2}\beta_u \sum_{h=1}^H \sqrt{\min(|I_h|, |J_h|)} \cdot \left(\sum_{k=1}^K \sum_{(i,j) \in X_h^k} \|\Phi(C_h^k, e_h^k, i, j)\|_{(\Sigma_h^{k+1})^{-1}} \right), \end{aligned} \quad (25)$$

where the second step is by $\min(|I_h|, |J_h|) \geq 1$. To bound the summation in the bracket, we seek to use the elliptical potential lemma. However, note that the telescoping sum involves several utility observations (i.e. all $(i, j) \in X_h^k$) at one time, instead of a single observation. To address this issue, for any (k, h) , let's index all the pairs $(i, j) \in X_h^k$ by an index $m = 1, \dots, |X_h^k|$. Here the order of the indexing does not matter. With a slight abuse of notation, we denote $\Phi_h^k(m) = \Phi(C_h^k, e_h^k, i, j)$ if (i, j) has index m in our indexing. We also define, for $n = 1, \dots, |X_h^k|$, that

$$\Sigma_h^{k+1}(n) = \Sigma_h^k + \sum_{m=1}^n \Phi_h^k(m) \Phi_h^k(m)^\top.$$

By the above definition, we have $\Sigma_h^k(n_1) \preceq \Sigma_h^k(n_2)$ for all $1 \leq n_1 < n_2 \leq |X_h^k|$, and $\Sigma_h^{k+1} = \Sigma_h^{k+1}(|X_h^k|)$. It follows that

$$\begin{aligned} \sum_{k=1}^K \sum_{(i,j) \in X_h^k} \|\Phi(C_h^k, e_h^k, i, j)\|_{(\Sigma_h^{k+1})^{-1}} &= \sum_{k=1}^K \sum_{(i,j) \in X_h^k} \|\Phi(C_h^k, e_h^k, i, j)\|_{(\Sigma_h^{k+1}(|X_h^k|))^{-1}} \\ &= \sum_{k=1}^K \sum_{m=1}^{|X_h^k|} \|\Phi_h^k(m)\|_{(\Sigma_h^{k+1}(|X_h^k|))^{-1}} \\ &\leq \sum_{k=1}^K \sum_{m=1}^{|X_h^k|} \|\Phi_h^k(m)\|_{(\Sigma_h^{k+1}(m))^{-1}} \\ &\leq \sqrt{K \cdot \min(|I_h|, |J_h|) \cdot d^2 \log \left(\frac{K \min(|I_h|, |J_h|) + d^2}{d^2} \right)}, \end{aligned}$$

where the first inequality is by $\Sigma_h^k(|X_h^k|)^{-1} \preceq \Sigma_h^k(m)^{-1}$, and the second inequality is by Lemma G.2 with $\lambda = 1$, $\Phi \in \mathbb{R}^{d^2}$, and $|X_h^k| \leq \min(|I_h|, |J_h|)$. Combining with (25), we get that

$$\sum_{k=1}^K \tilde{R}_k^M \leq 4\sqrt{2}\beta_u d \left(\sum_{h=1}^H \min(|I_h|, |J_h|) \right) \cdot \sqrt{K \log \left(\frac{K \min(|I_h|, |J_h|) + d^2}{d^2} \right)}.$$

The final bound follows by plugging in the expression of β_u , and using the fact that $W_h \leq \min\{|I_h|, |J_h|\}$. The $1 - 2\delta$ probability comes from the union bound of the event of Lemma C.1 and the event of (23). \square

E Proof of Lemmas

E.1 Proof of Lemma C.1

Proof of Lemma C.1. Fix arbitrary h . Since by Assumption 5.2, the noises in the observed utilities are independent and 1-sub-Gaussian, we can apply Theorem 2 in (Abbasi-Yadkori et al., 2011). We then have that, with probability at least $1 - \delta/2$, for any $k \in [K]$,

$$\|\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h\|_{\Sigma_h^k} \leq \sqrt{d^2 \cdot \log\left(\frac{1 + k \cdot \min\{|I_h|, |J_h|\} \cdot d^2/\lambda}{\delta/2}\right)} + \sqrt{\lambda}d.$$

We then have that for any C, e, i, j ,

$$\begin{aligned} u_h^k(C, e, i, j) - u_h(C, e, i, j) &= \langle \Phi(C, e, i, j), \boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h \rangle + \beta_u \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}} \\ &= \langle (\Sigma_h^k)^{-1/2} \Phi(C, e, i, j), (\Sigma_h^k)^{1/2} (\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h) \rangle + \beta_u \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}} \\ &\geq \beta_u \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}} - \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}} \cdot \|\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h\|_{\Sigma_h^k} \\ &\geq 0, \end{aligned}$$

where the first inequality follows from the Cauchy-Schwarz inequality and the second in equality is my the choice of β_u . Similarly, we have

$$\begin{aligned} u_h^k(C, e, i, j) - u_h(C, e, i, j) &\leq \beta_u \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}} + \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}} \cdot \|\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h\|_{\Sigma_h^k} \\ &\leq 2\beta_u \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}}. \end{aligned}$$

The same argument holds for $v_h^k - v_h$ with probability at least $1 - \delta/2$.

Finally, we take a union bound and conclude that the event holds with probability at least $1 - \delta$. \square

E.2 Proof of Lemma C.2

Proof of Lemma C.2. Note that \bar{r}_h is the maximum value of the linear program (4) with (I_h, J_h, u_h, v_h) , and \bar{r}_h^k is the maximum value of (4) with (I_h, J_h, u_h^k, v_h^k) . Since $u_h \leq u_h^k$ and $v_h \leq v_h^k$ by Lemma C.1 and the weights in (4) are restricted to be, it immediately holds that $\bar{r}_h(C, e, I_h, J_h) \leq \bar{r}_h^k(C, e, I_h, J_h)$.

On the other hand, denote by X_h the matching corresponding to \bar{r}_h . It follows that

$$\begin{aligned} \bar{r}_h(C, e, I_h, J_h) &= \sum_{(i,j) \in X_h} (u_h(C, e, i, j) + v_h(C, e, i, j)) \\ &\geq \sum_{(i,j) \in X_h^k} (u_h(C, e, i, j) + v_h(C, e, i, j)) \\ &= \sum_{(i,j) \in X_h^k} (u_h^k(C, e, i, j) + v_h^k(C, e, i, j)) \\ &\quad - \sum_{(i,j) \in X_h^k} (b_{u,h}(C, e, i, j) + b_{v,h}(C, e, i, j)), \end{aligned}$$

where the inequality is due to the sub-optimality of X_h^k under (u_h, v_h) . The result follows by using

$$\bar{r}_h^k(C, e, I_h, J_h) = \sum_{(i,j) \in X_h^k} (u_h^k(C, e, i, j) + v_h^k(C, e, i, j)),$$

and rearranging the terms. This completes the proof. \square

E.3 Proof of Lemma C.3

Lemma C.3 is a restatement of Lemma 5.4 in (Jagadeesan et al., 2021). For completeness, we present the proof here. Specifically, we prove a general version of Lemma C.3, which is Lemma E.1 below.

Lemma E.1. Let (u, v) and (\hat{u}, \hat{v}) be two pairs of utility functions on the agents set I, J , such that each of u, v, \hat{u}, \hat{v} maps from $I \times J$ to \mathbb{R} . Let $(\hat{X}, \hat{\tau})$ be a stable matching on (I, J) w.r.t. the utility functions (\hat{u}, \hat{v}) . Suppose $u \leq \hat{u}, v \leq \hat{v}$. Then the Subset Instability of $(\hat{X}, \hat{\tau})$ w.r.t. the utility (u, v) satisfies

$$\text{SI}(\hat{X}, \hat{\tau}; I, J, u, v) \leq \sum_{i \in I} \left| \hat{u}(i, \hat{X}(i)) - u(i, \hat{X}(i)) \right| + \sum_{j \in J} \left| \hat{v}(\hat{X}(j), j) - v(\hat{X}(j), j) \right|.$$

Proof of Lemma E.1 Define the function

$$\begin{aligned} f(I', J', X, \tau; u, v) &= \max_{X'} \left(\sum_{i \in I'} u(i, X'(i)) + \sum_{j \in J'} v(X'(j), j) \right) \\ &\quad - \sum_{i \in I'} (u(i, X(i)) + \tau(i)) - \sum_{j \in J'} (v(X(j), j) + \tau(j)), \end{aligned}$$

where $I' \times J' \subset I \times J$. Then by Definition 4.1 of Subset Instability, $\text{SI}(X, \tau; I, J, u, v) = \max_{I' \times J' \subset I \times J} f(I', J', X, \tau; u, v)$, and thus

$$\text{SI}(\hat{X}, \hat{\tau}; I, J, u, v) - \text{SI}(\hat{X}, \hat{\tau}; I, J, \hat{u}, \hat{v}) \leq \max_{I' \times J' \subset I \times J} \left[f(I', J', \hat{X}, \hat{\tau}; u, v) - f(I', J', \hat{X}, \hat{\tau}; \hat{u}, \hat{v}) \right].$$

To bound $f(I', J', \hat{X}, \hat{\tau}; u, v) - f(I', J', \hat{X}, \hat{\tau}; \hat{u}, \hat{v})$, we decompose

$$\begin{aligned} &f(I', J', \hat{X}, \hat{\tau}; u, v) - f(I', J', \hat{X}, \hat{\tau}; \hat{u}, \hat{v}) \\ &= \max_{X'} \left(\sum_{i \in I'} u(i, X'(i)) + \sum_{j \in J'} v(X'(j), j) \right) - \max_{X'} \left(\sum_{i \in I'} \hat{u}(i, X'(i)) + \sum_{j \in J'} \hat{v}(X'(j), j) \right) \\ &\quad + \sum_{i \in I'} \left(\hat{u}(i, \hat{X}(i)) + \hat{\tau}(i) \right) + \sum_{j \in J'} \left(\hat{v}(\hat{X}(j), j) + \hat{\tau}(j) \right) - \sum_{i \in I'} \left(u(i, \hat{X}(i)) + \hat{\tau}(i) \right) \\ &\quad - \sum_{j \in J'} \left(v(\hat{X}(j), j) + \hat{\tau}(j) \right) \\ &= \underbrace{\max_{X'} \left(\sum_{i \in I'} u(i, X'(i)) + \sum_{j \in J'} v(X'(j), j) \right) - \max_{X'} \left(\sum_{i \in I'} \hat{u}(i, X'(i)) + \sum_{j \in J'} \hat{v}(X'(j), j) \right)}_{\text{I}} \\ &\quad + \underbrace{\sum_{i \in I'} \left(\hat{u}(i, \hat{X}(i)) - u(i, \hat{X}(i)) \right) + \sum_{j \in J'} \left(\hat{v}(\hat{X}(j), j) - v(\hat{X}(j), j) \right)}_{\text{II}} \end{aligned}$$

Term I is nonpositive. To see this, note that by assumption $u \leq \hat{u}$ and $v \leq \hat{v}$. Thus the max-weight matching on $I' \times J'$ w.r.t. the utility functions (u, v) cannot exceed the max-weight matching on $I' \times J'$ w.r.t. (\hat{u}, \hat{v}) .

To bound term II, note that all the transfers $\hat{\tau}(i)$ and $\hat{\tau}(j)$ in the expression cancel out, and it follows that

$$\text{II} \leq \sum_{i \in I} \left(\hat{u}(i, \hat{X}(i)) - u(i, \hat{X}(i)) \right) + \sum_{j \in J} \left(\hat{v}(\hat{X}(j), j) - v(\hat{X}(j), j) \right)$$

where the inequality follows from the assumption that $u \leq \hat{u}$ and $v \leq \hat{v}$ and $I' \times J' \subseteq I \times J$. This finishes the proof. \square

Proof of Lemma C.3 For any fixed $k \in [K]$ and $h \in [H]$, replace I, J in Lemma E.1 with I_h, J_h , and replace u, v, \hat{u}, \hat{v} in with $u_h(C_h^k, e_h^k, \cdot, \cdot), v_h(C_h^k, e_h^k, \cdot, \cdot), u_h^k(C_h^k, e_h^k, \cdot, \cdot), v_h^k(C_h^k, e_h^k, \cdot, \cdot)$. Since X_h^k is the stable matching w.r.t. $u_h^k(C_h^k, e_h^k, \cdot, \cdot), v_h^k(C_h^k, e_h^k, \cdot, \cdot)$, we can replace \hat{X} with X_h^k . It then

follows from Lemma E.1 that

$$\begin{aligned}
\text{SI}_h^k &\leq \sum_{i \in I_h} |u_h^k(C_h^k, e_h^k, i, X_h^k(i)) - u_h(C_h^k, e_h^k, i, X_h^k(i))| \\
&\quad + \sum_{j \in J_h} |v_h^k(C_h^k, e_h^k, X_h^k(j), j) - v_h(C_h^k, e_h^k, X_h^k(j), j)| \\
&= \sum_{(i,j) \in X_h^k} |u_h^k(C_h^k, e_h^k, i, j) - u_h(C_h^k, e_h^k, i, j)| + |v_h^k(C_h^k, e_h^k, i, j) - v_h(C_h^k, e_h^k, i, j)| \\
&= \sum_{(i,j) \in X_h^k} (b_{u,h}(C_h^k, e_h^k, i, j) + b_{v,h}(C_h^k, e_h^k, i, j)),
\end{aligned}$$

where the second step holds because the true utility and the estimated utility are zero for unmatched agents under X_h^k , and the last step is by the definition of the bonus function $b_{u,h}$ and $b_{v,h}$. \square

E.4 Proof of Planner's Regret Decomposition

We first restate the lemma in its complete form.

Lemma C.5. *The planner's regret defined by (8) can be decomposed as*

$$\begin{aligned}
R^P(K) &= \underbrace{\sum_{k=1}^K \sum_{h=1}^H [\mathbb{E}_{\pi^*} [\delta_h^k(C_h, I_h, J_h, e_h) \mid C_1 = C_1^k] - \delta_h^k(C_h^k, I_h, J_h, e_h^k)]}_{E_1} + \underbrace{\sum_{k=1}^K \sum_{h=1}^H (\zeta_{k,h}^1 + \zeta_{k,h}^2)}_{E_2} \\
&\quad + \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{\pi^*} [\langle \bar{Q}_h^k(C_h, I_h, J_h, \cdot), \pi_h^*(\cdot \mid C_h, I_h, J_h) - \pi_{k,h}(\cdot \mid C_h, I_h, J_h) \rangle_{\Upsilon} \mid C_1 = C_1^k]}_{E_3},
\end{aligned}$$

where the expectation is over the trajectory $\{C_h, e_h\}_{h \in [H]}$ induced by executing the policy π^* (on the choice of $e \in \Upsilon$ only), and conditioning on $\{I_h, J_h\}_{h \in [H]}$ being fixed.

Proof of Lemma C.5. Recall the definition of the planner's regret from (8). We write

$$\bar{V}_1^*(s_1^k) - \bar{V}_1^{\pi^k}(s_1^k) = \underbrace{\bar{V}_1^*(s_1^k) - \bar{V}_1^k(s_1^k)}_{\text{I}} + \underbrace{\bar{V}_1^k(s_1^k) - \bar{V}_1^{\pi^k}(s_1^k)}_{\text{II}},$$

where $s_h^k = (C_h^k, I_h, J_h)$ by our notation.

Term I. We define two operators \mathbb{J}_h^* and $\mathbb{J}_{k,h}$ as

$$\mathbb{J}_h^* Q(s) = \langle Q(s, \cdot), \pi_h^*(\cdot \mid s) \rangle_{\Upsilon}, \quad \mathbb{J}_{k,h}^* Q(s) = \langle Q(s, \cdot), \pi_{k,h}(\cdot \mid s) \rangle_{\Upsilon},$$

for all $(k, h) \in [H] \times [K]$, $s \in \mathcal{C} \times 2^{\mathcal{I}} \times 2^{\mathcal{J}}$, and function $Q : \mathcal{C} \times \mathcal{I} \times \mathcal{J} \times \Upsilon \rightarrow \mathbb{R}$. Then by definition, we have $\bar{V}_h^k = \mathbb{J}_{k,h} \bar{Q}_h^k$, and $\bar{V}_h^* = \mathbb{J}_h^* \bar{Q}_h^*$. It follows that

$$\bar{V}_h^* - \bar{V}_h^k = \mathbb{J}_h^* \bar{Q}_h^* - \mathbb{J}_{k,h} \bar{Q}_h^k = \left(\mathbb{J}_h^* \bar{Q}_h^* - \mathbb{J}_h^* \bar{Q}_h^k \right) + \left(\mathbb{J}_h^* \bar{Q}_h^k - \mathbb{J}_{k,h} \bar{Q}_h^k \right) = \left(\mathbb{J}_h^* \bar{Q}_h^* - \mathbb{J}_h^* \bar{Q}_h^k \right) + \xi_h^k, \quad (26)$$

where $\xi_h^k := \mathbb{J}_h^* \bar{Q}_h^k - \mathbb{J}_{k,h} \bar{Q}_h^k$. Also, by the definition of δ_h^k in (18), we have

$$\bar{Q}_h^* - \bar{Q}_h^k = \bar{r}_h + \mathbb{P}_h \bar{V}_{h+1}^* - \left(\bar{r}_h + \mathbb{P}_h \bar{V}_{h+1}^k \right) + \delta_h^k = \mathbb{P}_h \left(\bar{V}_{h+1}^* - \bar{V}_{h+1}^k \right) + \delta_h^k.$$

Combining with (26), we get

$$\bar{V}_h^* - \bar{V}_h^k = \mathbb{J}_h^* \mathbb{P}_h \left(\bar{V}_{h+1}^* - \bar{V}_{h+1}^k \right) + \mathbb{J}_h^* \delta_h^k + \xi_h^k.$$

Applying the above equation recursively, we have that for any C_1, I_1, J_1 ,

$$\begin{aligned}
& \bar{V}_1^*(C_1, I_1, J_1) - \bar{V}_1^k(C_1, I_1, J_1) \\
&= \prod_{h=1}^H (\mathbb{J}_h^* \mathbb{P}_h) (\bar{V}_{H+1}^* - \bar{V}_{H+1}^k)(C_1, I_1, J_1) + \sum_{h=1}^H \left(\prod_{l=1}^{h-1} \mathbb{J}_l^* \mathbb{P}_l \right) \mathbb{J}_h^* \delta_h^k(C_1, I_1, J_1) \\
&\quad + \sum_{h=1}^H \left(\prod_{l=1}^{h-1} \mathbb{J}_l^* \mathbb{P}_l \right) \xi_h^k(C_1, I_1, J_1) \\
&= \mathbb{E}_{\pi^*} \left[\sum_{h=1}^H \delta_h^k(C_h, I_h, J_h, e_h) \middle| C_1, I_1, J_1 \right] \\
&\quad + \mathbb{E}_{\pi^*} \left[\sum_{h=1}^H \langle \bar{Q}_h^k(C_h, I_h, J_h, \cdot), [\pi_h^* - \pi_{k,h}] (\cdot \mid C_h, I_h, J_h) \rangle_{\Upsilon} \middle| C_1, I_1, J_1 \right], \quad (27)
\end{aligned}$$

where the second step holds because $\bar{V}_{H+1}^* = \bar{V}_{H+1}^k = 0$. By definition of the operators, it is clear that here the expectation \mathbb{E}_{π^*} is over the trajectory $\{(C_h, e_h)\}_{h \in [H]}$ induced by the planner executing the policy π^k to choose actions in Υ .

Term II. First note that by (18), the function $\delta_h^k : \mathcal{C} \times 2^{\mathcal{I}} \times 2^{\mathcal{J}} \times \Upsilon \rightarrow \mathbb{R}$ can be written as

$$\delta_h^k = \bar{r}_h + \mathbb{P}_h \bar{V}_{h+1}^k - \bar{Q}_h^k = \bar{r}_h + \mathbb{P}_h \bar{V}_{h+1}^k - \bar{Q}_h^{\pi_k} + \bar{Q}_h^{\pi_k} - \bar{Q}_h^k = \mathbb{P}_h (\bar{V}_{h+1}^k - \bar{V}_h^{\pi_k}) + (\bar{Q}_h^{\pi_k} - \bar{Q}_h^k), \quad (28)$$

where the last step is by $\bar{Q}_h^{\pi_k} = \bar{r}_h + \mathbb{P}_h \bar{V}_{h+1}^{\pi_k}$. Then for any h , we can write

$$\begin{aligned}
& [\bar{V}_h^k - \bar{V}_h^{\pi_k}] (C_h^k, I_h, J_h) \\
&= [\bar{V}_h^k - \bar{V}_h^{\pi_k} + \delta_h^k - \delta_h^k] (C_h^k, I_h, J_h) \\
&= [\bar{V}_h^k - \bar{V}_h^{\pi_k}] (C_h^k, I_h, J_h) + [\bar{Q}_h^{\pi_k} - \bar{Q}_h^k] (C_h^k, I_h, J_h, e_h^k) \\
&\quad + \mathbb{P}_h [\bar{V}_{h+1}^k - \bar{V}_{h+1}^{\pi_k}] (C_h^k, I_h, J_h, e_h^k) - \delta_h^k(C_h^k, I_h, J_h, e_h^k) \\
&= [\bar{V}_h^k - \bar{V}_h^{\pi_k}] (C_h^k, I_h, J_h) - [\bar{Q}_h^k - \bar{Q}_h^{\pi_k}] (C_h^k, I_h, J_h, e_h^k) \\
&\quad + \mathbb{P}_h [\bar{V}_{h+1}^k - \bar{V}_{h+1}^{\pi_k}] (C_h^k, I_h, J_h, e_h^k) - [\bar{V}_{h+1}^k - \bar{V}_{h+1}^{\pi_k}] (C_{h+1}^k, I_{h+1}, J_{h+1}) \\
&\quad + [\bar{V}_{h+1}^k - \bar{V}_{h+1}^{\pi_k}] (C_{h+1}^k, I_{h+1}, J_{h+1}) - \delta_h^k(C_h^k, I_h, J_h, e_h^k) \\
&= [\bar{V}_{h+1}^k - \bar{V}_{h+1}^{\pi_k}] (C_{h+1}^k, I_{h+1}, J_{h+1}) - \delta_h^k(C_h^k, I_h, J_h, e_h^k) \\
&\quad + \underbrace{[\mathbb{P}_h [\bar{V}_h^k - \bar{V}_h^{\pi_k}] (C_h^k, I_h, J_h, e_h^k) - [\bar{V}_{h+1}^k - \bar{V}_{h+1}^{\pi_k}] (C_{h+1}^k, I_{h+1}, J_{h+1})]}_{\zeta_{k,h}^2} \\
&\quad + \underbrace{[\bar{V}_h^k - \bar{V}_h^{\pi_k}] (C_h^k, I_h, J_h) - [\bar{Q}_h^k - \bar{Q}_h^{\pi_k}] (C_h^k, I_h, J_h, e_h^k)}_{\zeta_{k,h}^1},
\end{aligned}$$

where the second step is by (28). Applying the above equation recursively, we get

$$[\bar{V}_1^k - \bar{V}_1^{\pi_k}] (C_1^k, I_1, J_1) = \sum_{h=1}^H (\zeta_{k,h}^1 + \zeta_{k,h}^2) - \sum_{h=1}^H \delta_h^k(C_h^k, I_h, J_h, e_h^k), \quad (29)$$

where we use $\bar{V}_{H+1}^* = \bar{V}_{H+1}^k = 0$ again.

Combining (27) and (29), we get

$$\begin{aligned}
& \sum_{k=1}^K \left[\bar{V}_1^*(C_1^k, I_1, J_1) - \bar{V}_1^{\pi^k}(C_1^k, I_1, J_1) \right] \\
&= \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{\pi^*} \left[\delta_h^k(C_h, I_h, J_h, e_h) \middle| C_1, I_1, J_1 \right] - \sum_{k=1}^K \sum_{h=1}^H \delta_h^k(C_h^k, I_h, J_h, e_h^k) \\
&\quad + \sum_{k=1}^K \sum_{h=1}^H (\zeta_{k,h}^1 + \zeta_{k,h}^2) \\
&\quad + \mathbb{E}_{\pi^*} \left[\sum_{h=1}^H \langle \bar{Q}_h^k(C_h, I_h, J_h, \cdot), [\pi_h^* - \pi_{k,h}] \cdot \rangle \middle| C_1, I_1, J_1 \right],
\end{aligned}$$

which finishes the proof. \square

E.5 Proof of Lemma C.6

Proof of Lemma C.6 By the definition of \bar{Q}_h^k in Algorithm 1, the function δ_h^k satisfies

$$\begin{aligned}
\delta_h^k(C, I, J, e) &= \bar{r}_h(C, I, J, e) + \mathbb{P}_h \bar{V}_{h+1}^k(C, I, J, e) - \bar{Q}_h^k(C, I, J, e) \\
&= \bar{r}_h(C, I, J, e) + \mathbb{P}_h \bar{V}_{h+1}^k(C, I, J, e) - \bar{r}_h^k(C, I, J, e) - \hat{\mathbb{P}}_h \bar{V}_{h+1}^k(C, I, J, e).
\end{aligned}$$

In the sequel, we will show that \bar{r}_h^k and $\hat{\mathbb{P}}_h \bar{V}_{h+1}^k$ upper bound \bar{r}_h and $\mathbb{P}_h \bar{V}_{h+1}^k$ respectively.

We first consider the term $\bar{r}_h(C, e, I_h, J_h) - \bar{r}_h^k(C, e, I_h, J_h)$. It immediately follows from Lemma C.2 that, under the event of Lemma C.1

$$\begin{aligned}
-4 \sum_{(i,j) \in X_h^k} \beta_u \|\Phi(C, e, i, j)\|_{(\Sigma_h^k)^{-1}} &\leq - \sum_{(i,j) \in X_h^k} (b_{u,h}(C, e, i, j) + b_{v,h}(C, e, i, j)) \\
&\leq \bar{r}_h(C, e, I_h, J_h) - \bar{r}_h^k(C, e, I_h, J_h) \leq 0. \tag{30}
\end{aligned}$$

We now consider $\hat{\mathbb{P}}_h \bar{V}_{h+1}^k$. Since we are conditioning on $\{I_h, J_h\}$ which is independent of anything else, we can essentially treat it as a deterministic sequence. By (10), for any \bar{V}_{h+1}^k , there exists $\bar{\mathbf{w}}_h^k \in \mathbb{R}^d$ such that for any C, e , $\mathbb{P}_h \bar{V}_{h+1}^k(C, e, I_h, J_h) = \psi(C, e)^\top \bar{\mathbf{w}}_h^k$. This is because

$$\begin{aligned}
\mathbb{P}_h \bar{V}_{h+1}^k(C, e, I_h, J_h) &= \int \bar{V}_{h+1}^k(C', I_{h+1}, J_{h+1}) d\mathbb{P}(C' | C, e) \\
&= \int \bar{V}_{h+1}^k(C', I_{h+1}, J_{h+1}) \langle \psi(C, e), d\mu_h(C') \rangle \\
&= \langle \psi(C, e), \int \bar{V}_{h+1}^k(C', I_{h+1}, J_{h+1}) d\mu_h(C') \rangle.
\end{aligned}$$

Then we can write

$$\begin{aligned}
& \mathbb{P}_h \bar{V}_{h+1}^k(C, e, I_h, J_h) - \hat{\mathbb{P}}_h \bar{V}_{h+1}^k(C, e, I_h, J_h) \\
&= \psi(C, e)^\top \bar{\mathbf{w}}_h^k - \psi(C, e)^\top (\Lambda_h^k)^{-1} \sum_{t=1}^{k-1} \psi(C_h^t, e_h^t) \bar{V}_{h+1}^k(C_{h+1}^t, I_{h+1}, J_{h+1}) - \beta_V \cdot \|\psi(C, e)\|_{(\Lambda_h^k)^{-1}} \\
&= \psi(C, e)^\top (\Lambda_h^k)^{-1} \left[\Lambda_h^k \bar{\mathbf{w}}_h^k - \sum_{t=1}^{k-1} \psi(C_h^t, e_h^t) \bar{V}_{h+1}^k(C_{h+1}^t, I_{h+1}, J_{h+1}) \right] - \beta_V \cdot \|\psi(C, e)\|_{(\Lambda_h^k)^{-1}} \\
&= \psi(C, e)^\top (\Lambda_h^k)^{-1} \left[\sum_{t=1}^{k-1} \psi(C_h^t, e_h^t) \left(\mathbb{P}_h \bar{V}_{h+1}^k(C_h^t, e_h^t, I_h, J_h) - \bar{V}_{h+1}^k(C_{h+1}^t, I_{h+1}, J_{h+1}) \right) \right] \\
&\quad + \lambda \psi(C, e)^\top (\Lambda_h^k)^{-1} \bar{\mathbf{w}}_h^k - \beta_V \cdot \|\psi(C, e)\|_{(\Lambda_h^k)^{-1}}, \tag{31}
\end{aligned}$$

where the first step uses the construction of \mathbf{w}_h^k in Algorithm 1 and the last step uses the construction of Λ_h^k . For convenience, we write $s_h^t = (C_h^t, I_h, J_h)$ and $(s_h^t, e_h^t) = (C_h^t, e_h^t, I_h, J_h)$. It follows from the Cauchy-Schwarz inequality that the first part on the R.H.S. of (31) satisfies

$$\begin{aligned} & \left| \boldsymbol{\psi}(C, e)^\top (\Lambda_h^k)^{-1} \left[\sum_{t=1}^{k-1} \boldsymbol{\psi}(C_h^t, e_h^t) \left(\mathbb{P}_h \bar{V}_{h+1}^k(s_h^t, e_h^t) - \bar{V}_{h+1}^k(s_{h+1}^t) \right) \right] + \lambda \boldsymbol{\psi}(C, e)^\top (\Lambda_h^k)^{-1} \bar{\mathbf{w}}_h^k \right| \\ & \leq \|\boldsymbol{\psi}(C, e)\|_{(\Lambda_h^k)^{-1}} \cdot \left\| \sum_{t=1}^{k-1} \boldsymbol{\psi}(C_h^t, e_h^t) \left(\mathbb{P}_h \bar{V}_{h+1}^k(s_h^t, e_h^t) - \bar{V}_{h+1}^k(s_{h+1}^t) \right) \right\|_{(\Lambda_h^k)^{-1}} \\ & \quad + \lambda \|\boldsymbol{\psi}(C, e)\|_{(\Lambda_h^k)^{-1}} \cdot \|\bar{\mathbf{w}}_h^k\|_{(\Lambda_h^k)^{-1}}. \end{aligned} \quad (32)$$

In the following, to bound (31), we first bound the self-normalized stochastic process using tools from self-normalized martingale. The issue is that, according to Algorithm 1, the function \bar{V}_{h+1}^k depends on the first $(k-1)$ episodes and thus depends on the trajectory $\{(C_h^t, e_h^t, C_{h+1}^t)\}_{t \in [k-1]}$. We thus adopt a common approach to solve this issue by considering the function class containing each value function estimator \bar{V}_{h+1}^k . The covering trick is a commonly used technique (Ling et al., 2019), and we will discuss the detail of the construction of the function class and its covering in Section F. The covering trick allows us to get the following lemma.

Lemma E.2. *Under the setting of Theorem 5.5 with probability at least $1 - \delta$, for any $(h, k) \in [H] \times [K]$,*

$$\begin{aligned} & \left\| \sum_{t=1}^{k-1} \boldsymbol{\psi}(C_h^t, e_h^t) \left(\mathbb{P}_h \bar{V}_{h+1}^k(C_h^t, e_h^t, I_h, J_h) - \bar{V}_{h+1}^k(C_{h+1}^t, I_{h+1}, J_{h+1}) \right) \right\|_{(\Lambda_h^k)^{-1}} \\ & \leq 16d^2 \cdot \left(\sum_{h=1}^H W_h \right) \cdot \sqrt{\log \left(\frac{dKH \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot (\beta_V + \beta_u)}{\delta} \right)}. \end{aligned}$$

Proof of Lemma E.2 See Appendix E.7 for the proof. \square

For the term $\|\bar{\mathbf{w}}_h^k\|_{(\Lambda_h^k)^{-1}}$, by Assumption 5.2 and $|\bar{V}_h^k| \leq H \min\{|\mathcal{I}|, |\mathcal{J}|\}$, we have

$$\|\bar{\mathbf{w}}_h^k\|_2 = \left\| \int_{\mathcal{C}} \bar{V}_h^k(C, I_h, J_h) d\mu_h(C) \right\|_2 \leq \sqrt{d} \cdot \left(\sum_{h=1}^H W_h \right).$$

Combine the above inequality with Lemma E.2 and (32), and we get that

$$\begin{aligned} & \left| \boldsymbol{\psi}(C, e)^\top (\Lambda_h^k)^{-1} \left[\sum_{t=1}^{k-1} \boldsymbol{\psi}(C_h^t, e_h^t) \left(\mathbb{P}_h \bar{V}_{h+1}^k(s_h^t, e_h^t) - \bar{V}_{h+1}^k(s_{h+1}^t) \right) \right] + \lambda \boldsymbol{\psi}(C, e)^\top (\Lambda_h^k)^{-1} \bar{\mathbf{w}}_h^k \right| \\ & \leq 17d^2 \cdot \left(\sum_{h=1}^H W_h \right) \sqrt{\chi} \cdot \|\boldsymbol{\psi}(C, e)\|_{(\Lambda_h^k)^{-1}}, \end{aligned}$$

where

$$\chi := \log \left(\frac{dKH \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot (\beta_V + \beta_u)}{\delta} \right).$$

It remains to show that there exists choice of β_V (or equivalently, the constant η in the description of Theorem 5.5) such that

$$17d^2 \cdot \left(\sum_{h=1}^H W_h \right) \sqrt{\chi} \leq \beta_V.$$

Specifically, we show that we can pick some constant η and set

$$\beta_V = \eta d^2 \left(\sum_{h=1}^H W_h \right) \cdot \sqrt{\log \frac{dKH \min\{|\mathcal{I}|, |\mathcal{J}|\}}{\delta}}.$$

Indeed, plug in the expression of β_V and β_u and we get

$$\begin{aligned} \chi &\leq \log \left(\frac{3d^3 KH^2 \min\{|\mathcal{I}|, |\mathcal{J}|\}^2 \cdot \eta \log(3dKH \min\{|\mathcal{I}|, |\mathcal{J}|\}/\delta)}{\delta} \right) \\ &\leq 3 \log \left(\frac{3\eta dKH \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot \iota}{\delta} \right) \\ &= 3\iota + 3 \log(\eta) + 3 \log(\iota), \end{aligned}$$

where $\iota = \log(3dKH \min\{|\mathcal{I}|, |\mathcal{J}|\}/\delta)$. Since $\iota > \log 3$, it suffices to pick η such that

$$13\sqrt{3 \log 3 + 3 \log(\eta) + 3 \log(\log 3)} \leq \eta \cdot \log 3,$$

which finishes the proof. \square

E.6 Proof of Lemma C.7

Proof of Lemma C.7. The lemma can be proven by standard martingale concentration, similar to the analysis in (Cai et al., 2020; Yang et al., 2020). Specifically, we define the σ -fields as

$$\begin{aligned} \mathcal{F}_{k,h,0} &= \sigma \left(\{(C_l^t, e_l^t)_{(l,t) \in [k-1] \times [H]}\} \cup \{(C_l^k, e_l^k)\}_{l \in [H]} \right), \\ \mathcal{F}_{k,h,1} &= \sigma \left(\{(C_l^t, e_l^t)_{(l,t) \in [k-1] \times [H]}\} \cup \{(C_l^k, e_l^k)\}_{l \in [H]} \cup \{C_{h+1}^k\} \right). \end{aligned}$$

By definition, it is clear that these σ -fields form a filtration under the dictionary order on the index tuple (k, h, o) where $o \in \{0, 1\}$.

Note that for any $(k, h) \in [K] \times [H]$, since \bar{V}_h^k, \bar{Q}_h^k and the policy π_k are all functions of the first $(k-1)$ episodes, they are all $\mathcal{F}_{k,1,1}$ -measurable. As a result, $\zeta_{k,h}^1$ is $\mathcal{F}_{k,h,1}$ -measurable and $\zeta_{k,h}^2$ is $\mathcal{F}_{k,h,2}$ -measurable, for all (k, h) .

According to Algorithm 1, the planner's action $e_h^k \sim \pi_k(\cdot | C_h^k)$. This indicates that condition on C_h^k , $\bar{V}_h^{\pi_k}(C_h^k, I_h, J_h) - \bar{Q}_h^{\pi_k}(C_h^k, e_h^k, I_h, J_h) = 0$. Also, since $e_h^k \leftarrow \operatorname{argmax}_{e \in \Upsilon} \bar{Q}_h^k(s_h^k, e, I_h, J_h)$, and $\bar{V}_h^k(C, I, J) = \max_{e \in \Upsilon} \bar{Q}_h^k(C, e, I, J)$ for all (C, I, J) by the algorithm, we have $\bar{V}_h^k(C_h^k, I_h, J_h) - \bar{Q}_h^k(C_h^k, e_h^k, I_h, J_h) = 0$. Altogether we have $\zeta_{k,h}^1 = 0$ for all $(k, h) \in [K] \times [H]$.

For $\zeta_{k,h}^2$, first note that there is no dependence issue between the value functions and the trajectory since \bar{V}_{h+1}^k and $\bar{V}_{h+1}^{\pi_k}$ are functions of the first $(k-1)$ episodes. Then since $C_{h+1}^k \sim \mathbb{P}_h(\cdot | C_h^k, e_h^k)$, we have

$$\mathbb{E}[\zeta_{k,h}^2 | \mathcal{F}_{k,h,1}] = 0. \quad (33)$$

Thus we conclude that $\{(\zeta_{k,h}^1, \zeta_{k,h}^2)\}_{(k,h) \in [K] \times [H]}$ is a martingale difference sequence. Since $\bar{V}_h^k, \bar{Q}_h^k, \bar{V}_h^{\pi_k}, \bar{Q}_h^{\pi_k}$ are all bounded by W_h for all h, k , we apply the Azuma-Hoeffding inequality (Lemma G.3) and get that,

$$\mathbb{P}(|E_2| \geq \epsilon) \leq 2 \exp \left(\frac{-\epsilon^2}{8K \sum_{h=1}^H W_h^2} \right).$$

Equivalently, with probability at least $1 - \delta$, we have

$$|E_2| \leq \sqrt{8K \cdot \log \frac{2}{\delta}} \cdot \sqrt{\sum_{h=1}^H W_h^2} \leq \sqrt{8K \cdot \log \frac{2}{\delta}} \cdot \left(\sum_{h=1}^H W_h \right),$$

which finishes the proof. \square

E.7 Proof of Lemma E.2

Proof of Lemma E.2 By the analysis in Section F there exists a function class \mathcal{V} containing all \bar{V}_h^k , and the ϵ -covering number of \mathcal{V} is given by Lemma F.1. Also note that by the truncation, we have $|\bar{V}_h^k| \leq \sum_{l=h}^H W_l$. Then we apply Lemma G.5 with $R = \sum_{l=h}^H W_l$ and combine with Lemma F.1, and get that, fix any $0 < \epsilon < 1$, with probability at least $1 - \delta/H$, for all $k \in [K]$,

$$\begin{aligned} & \left\| \sum_{t=1}^{k-1} \psi(C_h^t, e_h^t) \left(\mathbb{P}_h \bar{V}_{h+1}^k(C_h^t, e_h^t, I_h, J_h) - \bar{V}_{h+1}^k(C_{h+1}^t, I_{h+1}, J_{h+1}) \right) \right\|_{(\Lambda_h^k)^{-1}}^2 \\ & \leq 4 \left(\sum_{l=h}^H W_l \right)^2 \left[\frac{d}{2} \log \left(\frac{k + \lambda}{\lambda} \right) + 6d^2 \log \left(1 + \frac{dkH \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot \beta_V}{\min\{\lambda, 1\} \cdot \epsilon} \right) \right. \\ & \quad \left. + 4d^4 \log \left(1 + \frac{d \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot \beta_u}{\min\{\lambda, 1\} \cdot \epsilon} \right) + \log \frac{1}{\delta} \right] + \frac{8k^2 \epsilon^2}{\lambda}. \end{aligned}$$

Let $\lambda = 1$, pick $\epsilon = d^2 \left(\sum_{l=h}^H W_l \right) / K$ and then take a union bound over $h \in [H]$, we get

$$\begin{aligned} & \left\| \sum_{t=1}^{k-1} \psi(C_h^t, e_h^t) \left(\mathbb{P}_h \bar{V}_{h+1}^k(C_h^t, e_h^t, I_h, J_h) - \bar{V}_{h+1}^k(C_{h+1}^t, I_{h+1}, J_{h+1}) \right) \right\|_{(\Lambda_h^k)^{-1}} \\ & \leq 16d^2 \cdot \left(\sum_{l=h}^H W_l \right) \cdot \sqrt{\log \left(\frac{dKH \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot (\beta_V + \beta_u)}{\delta} \right)}. \end{aligned}$$

□

F Covering Number of Function Classes

In this section, we will construct a function class \mathcal{V} that provably contains \bar{V}_h^k for all $(k, h) \in [K] \times [H]$. And we will compute the covering number of \mathcal{V} . The result is summarized by Lemma F.1 below.

Lemma F.1. *Assume $KH > 32$. For any $\epsilon < 1$, the ϵ -covering number of \mathcal{V} is upper bounded by*

$$\log \mathcal{N}_\epsilon^\mathcal{V} \leq 6d^2 \log \left(1 + \frac{dKH \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot \beta_V}{\min\{\lambda, 1\} \cdot \epsilon} \right) + 4d^4 \log \left(1 + \frac{d \min\{|\mathcal{I}|, |\mathcal{J}|\} \cdot \beta_u}{\min\{\lambda, 1\} \cdot \epsilon} \right).$$

To prove Lemma F.1, we will first construct a function class \mathcal{G} that contains all $\hat{\mathbb{P}}_h \bar{V}_{h+1}^k$, \mathcal{R} that contains all \bar{r}_h^k , and \mathcal{Q} that contains all \bar{Q}_h^k . The formal definition of these classes will be given in the following.

We also introduce the following technical lemma.

Lemma F.2 (Covering Number of ℓ_2 Ball). *For any $\epsilon > 0$, the ϵ -covering number of the ℓ_2 ball in \mathbb{R}^d with radius L is upper bound by $(1 + 2L/\epsilon)^d$.*

The proof of this classical result can be found in, for example, Chapter 5 in (Vershynin, 2010). Now we prove Lemma F.1

F.1 Proof of Lemma F.1

Covering of $\hat{\mathbb{P}}_h \bar{V}_{h+1}^k$. The next lemma is helpful to bound the covering number of the function class containing the function $\hat{\mathbb{P}}_h \bar{V}_{h+1}^k$. The proof is the same as that of Lemma D.6. in (Jin et al., 2020).

Lemma F.3. *Let $\mathcal{G} = \mathcal{G}(L, B)$ denote the function class with functions mapping from $\mathcal{C} \times \Upsilon$ to \mathbb{R} and of the following parametric form*

$$g(\cdot, \cdot) = \psi(\cdot, \cdot)^\top \mathbf{w} + \beta \cdot \|\psi(\cdot, \cdot)\|_{\Lambda^{-1}},$$

where $\|\mathbf{w}\| \leq L$, $\beta \in [0, B]$, and $\lambda_{\min}(\mathbf{\Lambda}) \geq \lambda > 0$. Assume $\|\psi(\cdot, \cdot)\|_2 \leq 1$. Let \mathcal{N}_ϵ denote the ϵ -covering number of \mathcal{G} with respect to the ℓ_∞ distance. Then we have

$$\log(\mathcal{N}_\epsilon^{\mathcal{G}}) \leq d \log(1 + 4L/\epsilon) + d^2 \log \left[1 + 8d^{1/2} B^2 / (\lambda \epsilon^2) \right].$$

Suppose for now that there exists $L_{\mathbf{w}} > 0$ such that $\|\mathbf{w}_h^k\|_2 \leq L_{\mathbf{w}}$ for all (k, h) . The value of $L_{\mathbf{w}}$ will be determined later. By applying Lemma F.3 with $L = L_{\mathbf{w}}$ and $B = \beta_V$, we get the following upper bound on the ϵ -covering number of the function class $\mathcal{G}(l_{\mathbf{w}}, \beta_V)$ which contains all $\widehat{\mathbb{P}}_h \bar{V}_{h+1}^k$:

$$\log(\mathcal{N}_\epsilon^{\mathcal{G}}) \leq d \log(1 + 4L_{\mathbf{w}}/\epsilon) + d^2 \log \left[1 + 8d^{1/2} \beta_V^2 / (\lambda \epsilon^2) \right]. \quad (34)$$

Covering of \bar{r}_h^k . We now define a function class \mathcal{R} which provably contains all the pseudo-reward estimates \bar{r}_h^k . Formally speaking, by Algorithm 2 the functions in \mathcal{R} are parametrized by the utility function estimates u_h^k and v_h^k . Denote the functions class containing all these utility function estimates by \mathcal{U} . Then according to Algorithm 4 any function $u : \mathcal{C} \times \Upsilon \times \mathcal{I} \times \mathcal{J} \rightarrow \mathbb{R}$ in \mathcal{U} can be written as

$$u(C, e, i, j) = \left(\langle \Phi(C, e, i, j), \theta \rangle + \beta_u \sqrt{\Phi(C, e, i, j)^\top \Sigma^{-1} \Phi(C, e, i, j)} \right)_{[-1,1]}. \quad (35)$$

where $\Phi \in \mathbb{R}^{d^2}$ satisfying $\|\Phi\|_1 \leq 1$ by Assumption 5.2, $\|\theta\| \leq L_u$ for some $L_u > 0$ to be determined, and $\Sigma \in \mathbb{R}^d \times \mathbb{R}^d$ such that $\lambda_{\min}(\Sigma) \geq \lambda$. Since the truncation is a contraction mapping, by Lemma F.3, the ϵ -covering number of \mathcal{U} is upper bounded by

$$\log(\mathcal{N}_\epsilon^{\mathcal{U}}) \leq d^2 \log(1 + 4L_u/\epsilon) + d^4 \log \left[1 + 8d\beta_u^2 / (\lambda \epsilon^2) \right]. \quad (36)$$

We now consider the function class \mathcal{R} . We formally define \mathcal{R} to be the function class such that any function $r \in \mathcal{R}$ can be represented by

$$r(C, e, I, J) = \text{RE}(u(C, e, \cdot, \cdot), v(C, e, \cdot, \cdot), I, J)$$

for some $u, v \in \mathcal{U}$. Let functions $r_1, r_2 \in \mathcal{R}$ be parametrized by u_1, v_1 and u_2, v_2 respectively, such that

$$\begin{aligned} r_1(C, e, I, J) &= \text{RE}(u_1(C, e, \cdot, \cdot), v_1(C, e, \cdot, \cdot), I, J), \\ r_2(C, e, I, J) &= \text{RE}(u_2(C, e, \cdot, \cdot), v_2(C, e, \cdot, \cdot), I, J), \end{aligned}$$

for all $C \in \mathcal{C}, e \in \mathcal{v}, I \subset \mathcal{I}$ and $J \subset \mathcal{J}$. According to the linear program 4 there exist some weights $w_1 = \{w_{1,i,j}\}_{(i,j) \in I \times J}$ and $w_2 = \{w_{2,i,j}\}_{(i,j) \in I \times J}$, such that

$$\begin{aligned} r_1(C, e, I, J) &= \sum_{(i,j) \in I \times J} w_{1,i,j} [u_1(C, e, i, j) + v_1(C, e, i, j)], \\ r_2(C, e, I, J) &= \sum_{(i,j) \in I \times J} w_{2,i,j} [u_2(C, e, i, j) + v_2(C, e, i, j)]. \end{aligned}$$

It follows that

$$\begin{aligned} (r_1 - r_2)(C, e, I, J) &= \sum_{(i,j) \in I \times J} w_{1,i,j} [u_1(C, e, i, j) + v_1(C, e, i, j)] \\ &\quad - \sum_{(i,j) \in I \times J} w_{2,i,j} [u_2(C, e, i, j) + v_2(C, e, i, j)] \\ &\leq \sum_{(i,j) \in I \times J} w_{2,i,j} [u_1(C, e, i, j) + v_1(C, e, i, j)] \\ &\quad - \sum_{(i,j) \in I \times J} w_{2,i,j} [u_2(C, e, i, j) + v_2(C, e, i, j)] \\ &\leq \sum_{(i,j) \in I \times J} w_{2,i,j} [(u_1(C, e, i, j) - u_2(C, e, i, j)) + (v_1(C, e, i, j) - v_2(C, e, i, j))] \\ &\leq \min\{|I|, |J|\} \cdot (\|u_1 - u_2\|_\infty + \|v_1 - v_2\|_\infty), \end{aligned}$$

where the second step holds because w_1 is the optimal weight given u_1 and v_1 and w_2 satisfies the constraint of the linear program with u_1 and v_1 . The same upper bound holds for the difference $(r_2 - r_1)$. Therefore, for any I, J , we have

$$\|r_1(\cdot, \cdot, I, J) - r_2(\cdot, \cdot, I, J)\|_\infty \leq \min\{|I|, |J|\} \cdot (\|u_1 - u_2\|_\infty + \|v_1 - v_2\|_\infty).$$

Since $I \subset \mathcal{I}$ and $J \subset \mathcal{J}$, in order for $\|r_1 - r_2\|_\infty \leq \epsilon$ to hold, it suffices to have $\|u_1 - u_2\|_\infty \leq \epsilon'$ and $\|v_1 - v_2\|_\infty \leq \epsilon'$ where $\epsilon' = \epsilon / (2 \min\{|\mathcal{I}|, |\mathcal{J}|\})$. Therefore, by (36), the ϵ -covering number of \mathcal{R} satisfies

$$\begin{aligned} \log \mathcal{N}_\epsilon^{\mathcal{R}} &\leq 2 \log \mathcal{N}_{\epsilon'}^u \leq 2d^2 \log(1 + 4L_u/\epsilon') + 2d^4 \log\left[1 + 8d\beta_u^2/(\lambda\epsilon'^2)\right] \\ &\leq 2d^2 \log\left(1 + \frac{8L_u \min\{|\mathcal{I}|, |\mathcal{J}|\}}{\epsilon}\right) + 2d^4 \log\left[1 + \frac{32d\beta_u^2 (\min\{|\mathcal{I}|, |\mathcal{J}|\})^2}{\lambda\epsilon^2}\right]. \end{aligned} \quad (37)$$

In the above we have shown that the function class \mathcal{R} contains all \bar{r}_h^k and \mathcal{G} contains all $\hat{\mathbb{P}}_h \bar{V}_{h+1}^k$. We now define the function class $\mathcal{Q} := \mathcal{R} + \mathcal{G}$ as

$$\mathcal{Q} := \left\{ (r + g)_{[0, \sum_{l=h}^H w_l]} \mid r \in \mathcal{R}, g \in \mathcal{G} \right\}.$$

Then it immediately follows from the algorithm that \mathcal{Q} contains all \bar{Q}_h^k functions. By (34) and (37), the ϵ -covering number of the function class \mathcal{Q} can be upper bounded by

$$\begin{aligned} \log \mathcal{N}_\epsilon^{\mathcal{Q}} &\leq d \log(1 + 4L_w/\epsilon) + d^2 \log\left[1 + 8d^{1/2}\beta_V^2/(\lambda\epsilon^2)\right] \\ &\quad + 2d^2 \log\left(1 + \frac{8L_u \min\{|\mathcal{I}|, |\mathcal{J}|\}}{\epsilon}\right) \\ &\quad + 2d^4 \log\left[1 + \frac{32d\beta_u^2 (\min\{|\mathcal{I}|, |\mathcal{J}|\})^2}{\lambda\epsilon^2}\right], \end{aligned} \quad (38)$$

Since by construction, $\bar{V}_h^k(C, I, J) = \max_e \bar{Q}_h^k(C, e, I, J)$ and taking the maximum is a contraction mapping, the upper bound in (38) also holds for $\log \mathcal{N}_\epsilon^{\mathcal{V}}$.

By Lemma D.1, we can pick

$$L_w = \left(\sum_{h=1}^H W_h \right) \cdot \sqrt{dK/\lambda} \quad \text{and} \quad L_u = \sqrt{\frac{d^2 K \cdot \min\{|\mathcal{I}|, |\mathcal{J}|\}}{\lambda}}.$$

From the above analysis, we can simplify the R.H.S. of (38) and get the desired bound for the covering number of \mathcal{V} . This finishes the proof of Lemma F.1.

G Auxiliary Lemmas

Lemma G.1 (Lemma D.1 in Jin et al. 2020). *For arbitrary d , let $\Lambda_k = \lambda \mathbf{I}_d + \sum_{t=1}^{k-1} \mathbf{x}_t \mathbf{x}_t^\top$ where $\mathbf{x}_t \in \mathbb{R}^d$ and $\lambda > 0$. Then*

$$\sum_{t=1}^{k-1} \mathbf{x}_t^\top (\Lambda_k)^{-1} \mathbf{x}_t \leq d.$$

Proof of Lemma G.1 We can write

$$\sum_{t=1}^{k-1} \mathbf{x}_t^\top (\Lambda_k)^{-1} \mathbf{x}_t = \sum_{t=1}^{k-1} \text{tr} \left(\mathbf{x}_t^\top (\Lambda_k)^{-1} \mathbf{x}_t \right) = \text{tr} \left((\Lambda_k)^{-1} \sum_{t=1}^{k-1} \mathbf{x}_t \mathbf{x}_t^\top \right).$$

Denote the eigenvalue of $\sum_{t=1}^{k-1} \mathbf{x}_t \mathbf{x}_t^\top$ as $\{\lambda_1, \dots, \lambda_d\}$, and decompose $\sum_{t=1}^{k-1} \mathbf{x}_t \mathbf{x}_t^\top = \mathbf{U} \text{diag}(\lambda_1, \dots, \lambda_d) \mathbf{U}^\top$. Then we have $\Lambda_k = \mathbf{U} \text{diag}(\lambda_1 + \lambda, \dots, \lambda_d + \lambda) \mathbf{U}^\top$. It follows that $\text{tr}((\Lambda_k)^{-1} \sum_{t=1}^{k-1} \mathbf{x}_t \mathbf{x}_t^\top) = \sum_{j=1}^d \lambda_j / (\lambda_j + \lambda) \leq d$. \square

The next is the well-known Elliptical Potential Lemma (Cesa-Bianchi and Lugosi, 2006; Abbasi-Yadkori et al., 2011; Lattimore and Szepesvári, 2020).

Lemma G.2 (Elliptical Potential Lemma). *For arbitrary d , let $\Lambda_k = \lambda \mathbf{I}_d + \sum_{t=1}^{k-1} \mathbf{x}_t \mathbf{x}_t^\top$ where $\mathbf{x}_t \in \mathbb{R}^d$ and $\lambda > 0$. Then*

$$\sum_{t=1}^k \|\mathbf{x}_t\|_{(\Lambda_{t+1})^{-1}} \leq \sqrt{kd \log \left(\frac{k+d\lambda}{d\lambda} \right)}.$$

Lemma G.3 (Azuma-Hoeffding inequality (Azuma, 1967)). *Let $\{X_t\}_{t=0}^\infty$ be a real-valued martingale such that for every $t \geq 1$, it holds that $|X_t - X_{t-1}| \leq B_t$ for some $B_t \geq 0$. Then*

$$\mathbb{P}(|X_t - X_0| \geq \epsilon) \leq 2 \exp \left(\frac{-\epsilon^2}{2 \sum_{\tau=1}^t B_\tau^2} \right).$$

G.1 Concentration Inequalities for Self-normalized Martingales

Theorem G.4 (Hoeffding inequality for Self-normalized martingales (Abbasi-Yadkori et al., 2011)). *Let $\{\eta_t\}_{t=1}^\infty$ be a real-valued stochastic process. Let $\{\mathcal{F}_t\}_{t=0}^\infty$ be a filtration, such that η_t is \mathcal{F}_t -measurable. Assume $\eta_t | \mathcal{F}_{t-1}$ is zero-mean and R -subgaussian for some $R > 0$, i.e.,*

$$\forall \lambda \in \mathbb{R}, \quad \mathbb{E} \left[e^{\lambda \eta_t | \mathcal{F}_{t-1}} \right] \leq e^{\lambda^2 R^2 / 2}.$$

Let $\{\mathbf{x}_t\}_{t=1}^\infty$ be an \mathbb{R}^d -valued stochastic process where \mathbf{x}_t is \mathcal{F}_{t-1} -measurable. Assume Λ_0 is a $d \times d$ positive definite matrix, and let $\Lambda_t = \Lambda_0 + \sum_{s=1}^t \mathbf{x}_s \mathbf{x}_s^\top$. Then, for any $\delta > 0$, with probability at least $1 - \delta$, for all $t > 0$,

$$\left\| \sum_{s=1}^t \mathbf{x}_s \eta_s \right\|_{\Lambda_t^{-1}}^2 \leq 2R^2 \log \left(\frac{\det(\Lambda_t)^{1/2} \det(\Lambda_0)^{-1/2}}{\delta} \right).$$

Lemma G.5 (Lemma D.4 in Jin et al., 2020). *Let \mathcal{V} be a function class such that any $V \in \mathcal{V}$ maps from $\mathcal{S} \rightarrow \mathbb{R}$ and $\|V\|_\infty \leq R$. Let $\{\mathcal{F}_t\}_{t=0}^\infty$ be a filtration. Let $\{s_t\}_{t=1}^\infty$ be a stochastic process in the space \mathcal{S} such that s_t is \mathcal{F}_t -measurable. Let $\{\mathbf{x}_t\}_{t=0}^\infty$ be an \mathbb{R}^d -valued stochastic process such that \mathbf{x}_t is \mathcal{F}_{t-1} -measurable and $\|\mathbf{x}_t\|_2 \leq 1$. Let $\Lambda_k = \lambda \mathbf{I} + \sum_{t=1}^{k-1} \mathbf{x}_t \mathbf{x}_t^\top$. Then for any $\delta > 0$, with probability at least $1 - \delta$, for any k , and any $V \in \mathcal{V}$, we have*

$$\left\| \sum_{t=1}^{k-1} \mathbf{x}_t [V(s_t) - \mathbb{E}[V(s_t) | \mathcal{F}_{t-1}]] \right\|_{(\Lambda_k)^{-1}}^2 \leq 4R^2 \left[\frac{d}{2} \log \left(\frac{k+\lambda}{\lambda} \right) + \log \frac{\mathcal{N}_\epsilon^\mathcal{V}}{\delta} \right] + \frac{8k^2 \epsilon^2}{\lambda},$$

where $\mathcal{N}_\epsilon^\mathcal{V}$ is the ϵ -covering number of \mathcal{V} with respect to the ℓ_∞ distance.

Proof of Lemma G.5 For any $V \in \mathcal{V}$, there exists V' in the ϵ -covering such that

$$V = V' + \Delta_V \quad \text{and} \quad \|\Delta_V\|_\infty \leq \epsilon.$$

Then we have

$$\begin{aligned} \left\| \sum_{t=1}^{k-1} \mathbf{x}_t [V(s_t) - \mathbb{E}[V(s_t) | \mathcal{F}_{t-1}]] \right\|_{(\Lambda_k)^{-1}}^2 &\leq 2 \left\| \sum_{t=1}^{k-1} \mathbf{x}_t [V'(s_t) - \mathbb{E}[V'(s_t) | \mathcal{F}_{t-1}]] \right\|_{(\Lambda_k)^{-1}}^2 \\ &\quad + 2 \left\| \sum_{t=1}^{k-1} \mathbf{x}_t [\Delta_V(s_t) - \mathbb{E}[\Delta_V(s_t) | \mathcal{F}_{t-1}]] \right\|_{(\Lambda_k)^{-1}}^2. \end{aligned}$$

For the first term on the R.H.S., we apply Theorem G.4 and a union bound to the ϵ -covering. The second term can be bound by $8k^2 \epsilon^2 / \lambda$ by using $\|\mathbf{x}_t\|_2 \leq 1$, $\lambda_{\min}(\Lambda_k) \geq \lambda$ and $\|\Delta_V\|_\infty \leq \epsilon$. \square