

Appendix of “Contextual Bandits with Knapsacks beyond Worst Cases via Re-Solving”

A From Discrete Randomness to Continuous Randomness

In the main body of this work, we explicitly assume that both the context set and the external factor set are discrete. Such an assumption can suitably capture most real-life situations. For example, in an agent’s online bidding problem with budget constraints, if we presume that the context is the agent’s actual value and the external factor is the highest competing bid, it is natural to suppose that all these three values are discrete. Nevertheless, for theoretical completeness, we expand our results in this section to circumstances where these two sets are infinite, i.e., the two underlying randomnesses is continuous. It is imperative to note that the scenario where one randomness is discrete and the other is continuous would be analogous in analysis by incorporating the techniques presented in Section 5.

Conceptually, the re-solving heuristic still works: we solve the optimization problem in each round concerning the remaining resources based on previous estimates. However, technically, since the distributions of context and external factors are continuous, we should further elaborate on the setting. In this section, we suppose that the context set $\Theta = [0, 1]^{d_u}$ and the external factor set $\Gamma = [0, 1]^{d_v}$. We denote $u(\theta)$ and $v(\gamma)$ as the density function of \mathcal{U} and \mathcal{V} , respectively. We assume that $p \in \{u, v\}$ belongs to the β_p -order L_p -Hölder smooth class $\Sigma(\beta_p, L_p)$. Here, for the foundation, given a vector $s = (s_1, \dots, s_d)$, define

$$|s| = s_1 + \dots + s_d, \quad D^s = \frac{\partial^{s_1 + \dots + s_d}}{\partial x_1^{s_1} \dots \partial x_d^{s_d}}.$$

Subsequently, for a positive integer β , the β -order L -Hölder smooth class is defined as

$$\Sigma(\beta, L) := \{g : |D^s g(x) - D^s g(y)| \leq L \|x - y\|_2, \text{ for all } s \text{ such that } |s| = \beta - 1, \text{ and all } x, y\}.$$

Now, suppose X_1, \dots, X_k are k i.i.d. samples from a distribution with density function $p \in \Sigma(\beta, L)$. According to Wasserman [2019], we have the following result, which implies that we can calculate an estimator from these samples that converges to the density function.

Proposition A.1 ([Wasserman, 2019]). *Suppose X_1, \dots, X_k are drawn i.i.d. from a d -dimension distribution \mathcal{P} , with density $p \in \Sigma(\beta, L)$ for some $L > 0$, and k is sufficiently large. Then there exists an estimator \hat{p}_k such that for any $\epsilon > 0$,*

$$\Pr \left[\sup_x |p(x) - \hat{p}_k(x)| > \frac{C \sqrt{\log(k/\epsilon)}}{k^{\beta/(2\beta+d)}} \right] \leq \epsilon,$$

with C a constant.

The details of constructing such a density estimator are postponed to Appendix F.1. We now return to the re-solving heuristic and Algorithm 1. In the algorithm, with continuous randomness, the constrained optimization problem to be solved in each round $\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)$ for $t = 1, 2, \dots$ becomes:

$$\begin{aligned} \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) &:= \max_{\phi: \Theta \times A^+ \rightarrow \mathbb{R}} \int_{\theta} \sum_{a \in A^+} \phi(\theta, a) \int_{\gamma} r(\theta, a, \gamma) \hat{v}_t(\gamma) \hat{u}_t(\theta) d\gamma d\theta, \\ \text{s.t. } &\int_{\theta} \sum_{a \in A^+} \phi(\theta, a) \int_{\gamma} c(\theta, a, \gamma) \hat{v}_t(\gamma) \hat{u}_t(\theta) d\gamma d\theta \leq \boldsymbol{\rho}_t, \\ &\sum_{a \in A^+} \phi(\theta, a) \leq 1, \quad \forall \theta \in \Theta, \\ &\phi(\theta, a) \geq 0, \quad \forall (\theta, a) \in \Theta \times A^+. \end{aligned}$$

Correspondingly, the reference optimization problem $J(\rho_t)$ is given below:

$$\begin{aligned}
J(\rho_t) &:= \max_{\phi: \Theta \times A^+ \rightarrow \mathbb{R}} \int_{\theta} \sum_{a \in A^+} \phi(\theta, a) \int_{\gamma} r(\theta, a, \gamma) v(\gamma) u(\theta) d\gamma d\theta, \\
\text{s.t. } &\int_{\theta} \sum_{a \in A^+} \phi(\theta, a) \int_{\gamma} c(\theta, a, \gamma) v(\gamma) u(\theta) d\gamma d\theta \leq \rho_t, \\
&\sum_{a \in A^+} \phi(\theta, a) \leq 1, \quad \forall \theta \in \Theta, \\
&\phi(\theta, a) \geq 0, \quad \forall (\theta, a) \in \Theta \times A^+.
\end{aligned}$$

At this point, it is worth mentioning that solving $\hat{J}(\rho_t, \mathcal{H}_t)$ in each round could be hard as it could be a continuous yet non-convex constrained optimization problem. Nevertheless, we assume the existence of an oracle that aids us in solving this optimization, and we focus on the regret of the re-solving method. Let $\alpha_u := (\beta_u + d_u)/(2\beta_u + d_u)$ and $\alpha_v := (\beta_v + d_v)/(2\beta_v + d_v)$, and we have the following two results, respectively, under full and partial information feedback.

Theorem A.1. *Under continuous randomness, with full information feedback, the expected accumulated reward Rew brought by Algorithm 1 satisfies:*

$$V^{\text{FL}} - \text{Rew} = O((T^{\alpha_u} + T^{\alpha_v} + T^{1/2})\sqrt{\log T}), \quad T \rightarrow \infty.$$

Theorem A.2. *Under continuous randomness, with partial information feedback, the expected accumulated reward Rew brought by Algorithm 1 satisfies:*

$$V^{\text{FL}} - \text{Rew} = O((T^{\alpha_u} + T^{1/2})\sqrt{\log T} + T^{\alpha_v} \log^{3/2-\alpha_v} T), \quad T \rightarrow \infty.$$

The proofs of the above theorems are presented in Appendices F.2 and F.3, respectively, which almost follow the threads of Theorems 5.1 and 5.2.

B Specifying the Worst-Case Location – Proof of Theorem 2.1

To prove the lemma, we first introduce an intermediate value, which we denote as V^{Hyb} , to upper bound V^{ON} , and show that the gap between V^{Hyb} and V^{FL} is $O(\sqrt{T})$ under the given condition. Specifically, we have the following definition:

$$\begin{aligned}
V^{\text{Hyb}} &:= \mathbb{E}_{\theta_1, \dots, \theta_T} \left[\max_{\phi_1, \dots, \phi_T: A^+ \rightarrow \mathbb{R}} \sum_{t=1}^T \sum_{a \in A^+} R(\theta_t, a) \phi_t(a) \right], \\
\text{s.t. } &\sum_{t=1}^T \sum_{a \in A^+} C(\theta_t, a) \phi_t(a) \leq \rho T, \\
&\sum_{a \in A^+} \phi_t(a) \leq 1, \quad \forall t \in [T], \\
&\phi_t(a) \geq 0, \quad \forall (t, a) \in [T] \times A^+.
\end{aligned} \tag{1}$$

To see that V^{Hyb} gives an upper bound on V^{ON} , we fix a request trajectory $\theta_1, \dots, \theta_T$. Now, for any non-anticipating strategy π , we let

$$p_t^{\pi}(a) = \Pr[a_t^{\pi} = a \mid \theta_1, \dots, \theta_t]$$

be the total probability that $a_t^{\pi} = a$ conditioning on the pre-determined request sequence, with respect to $\gamma_1, \dots, \gamma_{t-1}$ and the randomness of strategy π . We show that $\{p_t^{\pi}\}_{t=1, \dots, T}$ is a feasible solution to V^{Hyb} under $\theta_1, \dots, \theta_T$. Here, a key observation is that for any $t \in [T]$:

$$\begin{aligned}
\mathbb{E}[c(\theta_t, a_t^{\pi}, \gamma_t) \mid \theta_1, \dots, \theta_t] &= \mathbb{E}_{\gamma_t} \left[\sum_{a \in A^+} c(\theta_t, a, \gamma_t) \cdot \Pr[a_t^{\pi} = a \mid \theta_1, \dots, \theta_t] \right] \\
&= \sum_{a \in A^+} C(\theta_t, a) p_t^{\pi}(a).
\end{aligned}$$

460 In the above, the first expectation is taken on $\gamma_1, \dots, \gamma_t$ and the random choice of strategy π . Since
 461 $\sum_{t=1}^T \mathbf{c}(\theta_t, a_t^\pi, \gamma_t) \leq \boldsymbol{\rho}_T$ always holds, we derive that

$$\sum_{t=1}^T \sum_{a \in A^+} \mathbf{C}(\theta_t, a) p_t^\pi(a) = \mathbb{E} \left[\sum_{t=1}^T \mathbf{c}(\theta_t, a_t^\pi, \gamma_t) \mid \theta_1, \dots, \theta_T \right] \leq \boldsymbol{\rho}_T,$$

462 which indicates that $\{p_t^\pi\}_{t=1, \dots, T}$ is feasible to V^{Hyb} under $\theta_1, \dots, \theta_T$. To the same reason, we also
 463 have

$$\sum_{t=1}^T \sum_{a \in A^+} R(\theta_t, a) p_t^\pi(a) = \mathbb{E} \left[\sum_{t=1}^T \mathbf{r}(\theta_t, a_t^\pi, \gamma_t) \mid \theta_1, \dots, \theta_T \right]$$

464 equals the conditional expected reward of strategy π . Thus, since V^{Hyb} is a maximization problem
 465 for any request trajectory, we conclude that $V^{\text{Hyb}} \geq V^{\text{ON}}$.

466 It remains to show that when V^{FL} , or $J(\boldsymbol{\rho})$ has a unique and degenerate solution, $V^{\text{FL}} - V^{\text{Hyb}} =$
 467 $\Omega(\sqrt{T})$. We first present a transformation of V^{Hyb} . We let

$$x(\theta) := \frac{\#[\text{appearance of } \theta]}{T}$$

468 be the random variable indicating the frequency of θ when θ is drawn T times i.i.d. from \mathcal{U} . Obviously,
 469 the mean of \mathbf{x} is \mathbf{u} . We now demonstrate that

$$\begin{aligned} V^{\text{Hyb}} &= T \cdot \mathbb{E}_{\mathbf{x}} \left[\max_{\phi: \Theta \times A^+ \rightarrow \mathbb{R}} \sum_{\theta \in \Theta} x(\theta) \sum_{a \in A^+} R(\theta, a) \phi(\theta, a) \right], \\ \text{s.t. } &\sum_{\theta \in \Theta} x(\theta) \sum_{a \in A^+} \mathbf{C}(\theta, a) \phi(\theta, a) \leq \boldsymbol{\rho}, \\ &\sum_{a \in A^+} \phi(\theta, a) \leq 1, \quad \forall \theta \in \Theta, \\ &\phi(\theta, a) \geq 0, \quad \forall (\theta, a) \in \Theta \times A^+. \end{aligned} \quad (2)$$

470 To see this, in form (1), it is not hard to see that conditioning on $\theta_1, \dots, \theta_T$, the value of the
 471 optimization is only related to the number of times that any $\theta \in \Theta$ appears in the sequence, and
 472 irrelevant with their arriving order. Therefore, by taking an average, it is without loss of generality
 473 to suppose that $\phi_{t_1}^* = \phi_{t_2}^*$ as long as $\theta_{t_1} = \theta_{t_2}$. Under such an observation, it is natural that (1) is
 474 equivalent to (2).

475 For convenience, we now recall the definition of V^{FL} :

$$\begin{aligned} V^{\text{FL}} &= T \cdot \max_{\phi: \Theta \times A^+ \rightarrow \mathbb{R}} \sum_{\theta \in \Theta} u(\theta) \sum_{a \in A^+} R(\theta, a) \phi(\theta, a), \\ \text{s.t. } &\sum_{\theta \in \Theta} u(\theta) \sum_{a \in A^+} \mathbf{C}(\theta, a) \phi(\theta, a) \leq \boldsymbol{\rho}, \\ &\sum_{a \in A^+} \phi(\theta, a) \leq 1, \quad \forall \theta \in \Theta, \\ &\phi(\theta, a) \geq 0, \quad \forall (\theta, a) \in \Theta \times A^+. \end{aligned}$$

476 By Sierksma [2001], we know that when $J(\boldsymbol{\rho})$ has a unique and degenerate solution, then its dual
 477 form has multiple solutions. We then adopt the framework of Vera and Banerjee [2021]. In particular,
 478 we let $\boldsymbol{\lambda} \geq \mathbf{0}$ be the dual variable vector for the resource constraints, and $\boldsymbol{\mu} \geq \mathbf{0}$ be the dual
 479 variable vector for the probability feasibility constraints. If we take $\omega(\theta) = \mu(\theta)/u(\theta)$, then the dual
 480 programming of V^{FL}/T is the following as a function of \mathbf{u} :

$$\begin{aligned} \mathcal{D}[Z(\mathbf{u})] &= \min_{\boldsymbol{\lambda}, \boldsymbol{\omega}} \boldsymbol{\rho}^\top \boldsymbol{\lambda} + \mathbf{u}^\top \boldsymbol{\omega}, \\ \text{s.t. } &\boldsymbol{\lambda}^\top \mathbf{C}(\theta, a) + \omega(\theta) \geq R(\theta, a), \quad \forall (\theta, a) \in \Theta \times A^+, \\ &\boldsymbol{\lambda} \geq \mathbf{0}, \quad \boldsymbol{\omega} \geq \mathbf{0}. \end{aligned}$$

Now, suppose (λ^1, ω^1) and (λ^2, ω^2) are two different optimal solutions to $\mathcal{D}[Z(u)]$, which directly leads to $\lambda^1 \neq \lambda^2$ by the programming formation. We let $\lambda' = \lambda^1 - \lambda^2$ and $\omega' = \omega^1 - \omega^2$. Then,

$$\rho^\top \lambda^1 + u^\top \omega^1 = \rho^\top \lambda^2 + u^\top \omega^2 \implies \rho^\top \lambda' + u^\top \omega' = 0. \quad (3)$$

Further, notice that (λ^1, ω^1) and (λ^2, ω^2) are both feasible for $\mathcal{D}[Z(x)]$ for any x . Since $\mathcal{D}[Z(x)]$ is a minimization problem, by a convex combination, we have

$$\mathcal{D}[Z(x)] \leq (\rho^\top \lambda^1 + x^\top \omega^1) \mathbf{1}[\rho^\top \lambda' + x^\top \omega' \leq 0] + (\rho^\top \lambda^2 + x^\top \omega^2) \mathbf{1}[\rho^\top \lambda' + x^\top \omega' > 0].$$

Further, by optimality, we know that for any x ,

$$\mathcal{D}[Z(u)] = (\rho^\top \lambda^1 + u^\top \omega^1) \mathbf{1}[\rho^\top \lambda' + u^\top \omega' \leq 0] + (\rho^\top \lambda^2 + u^\top \omega^2) \mathbf{1}[\rho^\top \lambda' + u^\top \omega' > 0].$$

Now, by weak duality, since V^{Hyb}/T for any given x is a maximization problem, we know from the above two equations that

$$\begin{aligned} & (V^{\text{FL}} - V^{\text{Hyb}})/T \\ & \geq \mathcal{D}[Z(u)] - \mathbb{E}_x[\mathcal{D}[Z(x)]] \\ & \geq \mathbb{E}_x [((u - x)^\top \omega^1) \mathbf{1}[\rho^\top \lambda' + x^\top \omega' \leq 0] + ((u - x)^\top \omega^2) \mathbf{1}[\rho^\top \lambda' + x^\top \omega' > 0]] \\ & \stackrel{(a)}{=} \mathbb{E}_x [((u - x)^\top \omega^1) \mathbf{1}[(u - x)^\top \omega' \geq 0] + ((u - x)^\top \omega^2) (1 - \mathbf{1}[(u - x)^\top \omega' \geq 0])] \\ & \stackrel{(b)}{=} \mathbb{E}_x [((u - x)^\top \omega') \mathbf{1}[(u - x)^\top \omega' \geq 0]] \end{aligned}$$

Here, (a) is due to (3), and (b) is since the mean of x is u . Now, we let $\xi = \sqrt{T}(u - x)^\top \omega'$ be the normalized scaled variable. By Central Limit Theorem, $\xi \mathbf{1}[\xi \geq 0]$ converges to a half-normal distribution, which has constant expectation. Thus, we arrive at $V^{\text{FL}} - V^{\text{Hyb}} = \Omega(\sqrt{T})$, which finish the proof.

C Missing Proofs in Section 3

C.1 Proof of Theorem 3.1

We now give a proof of Theorem 3.1. The proof draws inspiration from that of Chen et al. [2022], but significantly diverges in terms of the problem setting.

C.1.1 Regret Decomposition

We start by presenting a regret decomposition approach, which stands on the dual viewpoint. We first recall the optimization problem $V^{\text{FL}} = T \cdot J(\rho_1)$:

$$\begin{aligned} J(\rho_1) &:= \max_{\phi: \Theta \times A^+ \rightarrow \mathbb{R}} \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} R(\theta, a) \phi(\theta, a) \right], \\ \text{s.t. } & \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} C(\theta, a) \phi(\theta, a) \right] \leq \rho_1, \\ & \sum_{a \in A^+} \phi(\theta, a) \leq 1, \quad \forall \theta \in \Theta, \\ & \phi(\theta, a) \geq 0, \quad \forall (\theta, a) \in \Theta \times A^+. \end{aligned}$$

Recall that $u(\theta)$ denotes the mass function of \mathcal{U} , then the above linear programming can be expanded as

$$\begin{aligned} J(\rho_1) &:= \max_{\phi: \Theta \times A^+ \rightarrow \mathbb{R}} \sum_{\theta \in \Theta, a \in A^+} u(\theta) R(\theta, a) \phi(\theta, a), \\ \text{s.t. } & \sum_{\theta \in \Theta, a \in A^+} u(\theta) C(\theta, a) \phi(\theta, a) \leq \rho_1, \\ & \sum_{a \in A^+} \phi(\theta, a) \leq 1, \quad \forall \theta \in \Theta, \\ & \phi(\theta, a) \geq 0, \quad \forall (\theta, a) \in \Theta \times A^+. \end{aligned}$$

Now let $\lambda \geq \mathbf{0}$ be the dual vector for the consumption constraint and $\{\mu^*(\theta)\}_{\theta \in \Theta} \geq \mathbf{0}$ be the dual variables for the action distribution constraint. By the strong duality of linear programming, there is an optimal dual variable tuple $(\lambda^*, \{\mu^*(\theta)\}_{\theta \in \Theta}) \geq \mathbf{0}$ such that:

$$\begin{aligned} J(\rho_1) &= \sum_{\theta \in \Theta, a \in A^+} (u(\theta) (R(\theta, a) - (\lambda^*)^\top C(\theta, a)) - \mu^*(\theta)) \phi_1^*(\theta, a) + (\lambda^*)^\top \rho_1 + \sum_{\theta \in \Theta} \mu^*(\theta) \\ &= \sum_{\theta \in \Theta, a \in A^+} u(\theta) (R(\theta, a) - (\lambda^*)^\top C(\theta, a)) \phi_1^*(\theta, a) + (\lambda^*)^\top \rho_1. \end{aligned} \quad (4)$$

Here ϕ_1^* is the optimal solution to $J(\rho_1)$. With (4), we have the following lemma for regret decomposition.

Lemma C.1. *For any stopping time $T_e \leq T_0$ adapted to the process $\{B_t\}$'s, we have*

$$\begin{aligned} &V^{\text{FL}} - \text{Rew} \\ &\leq \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} (u(\theta) (R(\theta, a) - (\lambda^*)^\top C(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \hat{\phi}_t^*(\theta, a)) \right] \\ &\quad + \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \right) \right] \\ &\quad + (\lambda^*)^\top \mathbb{E}[B_{T_e+1}] + \max_{\theta \in \Theta, a \in A^+} (R(\theta, a) - (\lambda^*)^\top C(\theta, a)) \cdot \mathbb{E}[T - T_e]. \end{aligned} \quad (5)$$

The proof of Lemma C.1 is deferred to Appendix C.2. We now give a brief explanation on this result. The first two terms in (5) depicts the gap between the choice of Algorithm 1 and the optimal decision. This is apparent for the first term. For the second term, we should notice that by complementary slackness, for each $\theta \in \Theta$,

$$\mu^*(\theta) \cdot \left(1 - \sum_{a \in A^+} \phi_1^*(\theta, a) \right) = 0.$$

Therefore, the second term in (5) is bounded if $\hat{\phi}_t^*$ is close to ϕ_1^* .

On the other hand, the last two terms are closely related to the choice of stopping time T_e and the consumption behavior of Algorithm 1. Intuitively, if T_e is sufficiently close to T , then $\mathbb{E}[T - T_e]$ should be appropriately bounded. Nevertheless, if the algorithm spends the resources too fast, then such a sufficiently large T_e would be impossible. Conversely, if the resources are consumed substantially slower than the optimal, then the term $\mathbb{E}[B_{T_e+1}]$, the remaining resources at the stopping time, would be unbounded.

In the following, we will deal with these two parts correspondingly. A crux to the analysis is to pick a satisfying stopping time T_e , which we will first cover.

C.1.2 The Gap to Optimal Decision

We first give a realization of the stopping time T_e , which relies on Assumption 3.1. As is shown by Mangasarian and Shiao [1987], Chen et al. [2022], local stability holds for an LP with unique and non-degenerate optimal solution, that is, the basic variables and binding constraints are kept within a minor perturbation on the coefficients. To this end, we first explicitly define the relevant concepts.

Definition C.1. A context-action pair (θ, a) is a basic variable for $J(\rho_1)$ if $\phi_1^*(\theta, a) > 0$, or else, it is a non-basic variable. Similarly define basic/non-basic variables for $\hat{J}(\rho_t, \mathcal{H}_t)$.

Definition C.2. $i \in [n]$ is a binding constraint for $J(\rho_1)$ if

$$\sum_{\theta \in \Theta, a \in A^+} u(\theta) C^i(\theta, a) \phi_1^*(\theta, a) = \rho_1^i,$$

or else it is a non-binding constraint. We let

$$\begin{aligned} \mathcal{S} &:= \{i \in [n] : i \text{ is a binding constraint for } J(\rho_1)\}, \\ \mathcal{T} &:= \{i \in [n] : i \text{ is a non-binding constraint for } J(\rho_1)\}, \end{aligned}$$

and use $\kappa|_{\mathcal{S}}$ to define the sub-vector of κ confined on \mathcal{S} , similar for $\kappa|_{\mathcal{T}}$. Further, $\theta \in \Theta$ is a binding constraint for $J(\rho_1)$ if

$$\sum_{a \in A^+} \phi_1^*(\theta, a) = 1,$$

or else it is a non-binding constraint. Similarly define binding/non-binding constraints for $\hat{J}(\rho_t, \mathcal{H}_t)$.

Under the above definitions, we have the following lemma, which is a derivation of the result in Chen et al. [2022]. We will provide the proof in Appendix C.3:

Lemma C.2 (Stability). *Under Assumption 3.1, there is a $D > 0$, such that when the following holds:*

$$\begin{aligned} \max \{ \|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_{\infty}, \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \} &\leq D, \\ \max \{ \|\rho_1|_{\mathcal{S}} - \rho_t|_{\mathcal{S}}\|_{\infty}, \max \{ \rho_1|_{\mathcal{T}} - \rho_t|_{\mathcal{T}} \} \} &\leq D, \end{aligned} \quad (6)$$

$J(\rho_1)$ and $\hat{J}(\rho_t, \mathcal{H}_t)$ share the same sets of basic/non-basic variables and binding/non-binding constraints.

With Lemma C.2 in hand, we can derive that when condition (6) is met, it holds that

$$(u(\theta) (R(\theta, a) - (\lambda^*)^{\top} C(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \hat{\phi}_t^*(\theta, a)) = 0, \quad (7)$$

$$\sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \right) = 0. \quad (8)$$

To see these, notice that by the dual feasibility of $J(\rho_1)$, we have $u(\theta) (R(\theta, a) - (\lambda^*)^{\top} C(\theta, a)) - \mu^*(\theta) \leq 0$. When $u(\theta) (R(\theta, a) - (\lambda^*)^{\top} C(\theta, a)) - \mu^*(\theta) < 0$, by primal optimality, $\phi_1^*(\theta, a) = 0$ and thus (θ, a) is non-basic for $J(\rho_1)$. By Lemma C.2, (θ, a) is also non-basic for $\hat{J}(\rho_t, \mathcal{H}_t)$ and $\hat{\phi}_t^*(\theta, a) = 0$ holds as well. This finishes the deduction of (7). A similar reasoning on binding constraints would help us achieve (8), which we omit here.

As the above goes, it is then natural for us to define T_e the stopping time in our analysis as follows:

$$T_e := \min\{T_0, \min\{t : \max\{\|\rho_1|_{\mathcal{S}} - \rho_t|_{\mathcal{S}}\|_{\infty}, \max\{\rho_1|_{\mathcal{T}} - \rho_t|_{\mathcal{T}}\}\} > D\} - 1\}, \quad (9)$$

where T_0 is the stopping time of Algorithm 1. That is to say, we always have $\max\{\|\rho_1|_{\mathcal{S}} - \rho_t|_{\mathcal{S}}\|_{\infty}, \max\{\rho_1|_{\mathcal{T}} - \rho_t|_{\mathcal{T}}\}\} \leq D$ when $t \leq T_e$. What we are left is to bound the situation when $\max\{\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_{\infty}, \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1\} > D$ for $1 \leq t \leq T_e$. In total, we arrive at the following result for this part, with the proof given in Appendix C.4:

Lemma C.3. *Under Assumption 3.1, with full information feedback, we have when $T \rightarrow \infty$:*

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} (u(\theta) (R(\theta, a) - (\lambda^*)^{\top} C(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \hat{\phi}_t^*(\theta, a)) \right] &= O(1), \\ \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \right) \right] &= O(1). \end{aligned}$$

We are now only left to bound the last two terms in (5).

C.1.3 The Gap to Optimal Consumption

As presented in (5), we now bound the remaining two terms, respectively $\mathbb{E}[B_{T_e+1}]$ and $\mathbb{E}[T - T_e]$ for T_e defined in (9). It turns out that these two terms are closely related. Due to this observation, we would first bound $(\lambda^*)^{\top} \cdot \mathbb{E}[B_{T_e+1}]$ by $\mathbb{E}[T - T_e]$, and then bound $\mathbb{E}[T - T_e]$.

Now by the strong duality of $J(\rho_1)$, we know that complementary slackness holds, that is $\lambda^*|_{\mathcal{T}} = \mathbf{0}$. We therefore have

$$\begin{aligned} (\lambda^*)^{\top} \mathbb{E}[B_{T_e+1}] &\leq (\lambda^*)^{\top} \mathbb{E}[B_{T_e}] = (\lambda^*|_{\mathcal{S}})^{\top} \mathbb{E}[B_{T_e}|_{\mathcal{S}}] = (\lambda^*|_{\mathcal{S}})^{\top} \mathbb{E}[(T - T_e + 1)\rho_{T_e}|_{\mathcal{S}}] \\ &\stackrel{(a)}{\leq} n(\rho^{\max} + D)\|\lambda^*\|_{\infty} \cdot \mathbb{E}[T - T_e + 1]. \end{aligned} \quad (10)$$

557 In the above, recall that ρ^{\max} denotes the maximum coordinate of $\boldsymbol{\rho}_1$, and D is specified in
 558 Lemma C.2. Consequently, (a) is due to the definition of T_e and that $\|\boldsymbol{\rho}_1 + D\mathbf{1}\|_\infty \leq \rho^{\max} + D$.

559 We are left to bound $\mathbb{E}[T - T_e]$. Nevertheless, this part would be rather technical and involved.
 560 Therefore we defer the analysis to Appendix C.5, and only give the final bounds.

561 **Lemma C.4.** *Under Assumption 3.1, with full information feedback, we have when $T \rightarrow \infty$:*

$$(\boldsymbol{\lambda}^*)^\top \mathbb{E}[\mathbf{B}_{T_e+1}] + \max_{\theta \in \Theta, a \in A^+} (R(\theta, a) - (\boldsymbol{\lambda}^*)^\top \cdot \mathbf{C}(\theta, a)) \mathbb{E}[T - T_e] = O(1).$$

562 Combining Lemmas C.1, C.3 and C.4, we arrive at Theorem 3.1.

563 C.2 Proof of Lemma C.1

564 The proof is obtained by the following set of (in)equalities.

$$\begin{aligned} & V^{\text{FL}} - \text{Rew} \\ &= T \cdot J(\boldsymbol{\rho}_1) - \mathbb{E} \left[\sum_{t=1}^{T_0} r(\theta_t, a_t, \gamma_t) \right] \\ &\stackrel{(a)}{\leq} T \cdot J(\boldsymbol{\rho}_1) - \mathbb{E} \left[\sum_{t=1}^{T_e} r(\theta_t, a_t, \gamma_t) \right] \\ &\stackrel{(b)}{=} T \cdot J(\boldsymbol{\rho}_1) - \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} u(\theta) R(\theta, a) \hat{\phi}_t^*(\theta, a) \right] \\ &\stackrel{(c)}{=} T \cdot \left(\sum_{\theta \in \Theta, a \in A^+} (u(\theta)(R(\theta, a) - (\boldsymbol{\lambda}^*)^\top \mathbf{C}(\theta, a)) - \mu^*(\theta)) \phi_1^*(\theta, a) + (\boldsymbol{\lambda}^*)^\top \boldsymbol{\rho}_1 + \sum_{\theta \in \Theta} \mu^*(\theta) \right) \\ &\quad - \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} u(\theta) R(\theta, a) \hat{\phi}_t^*(\theta, a) \right] \\ &\stackrel{(d)}{=} \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} (u(\theta)(R(\theta, a) - (\boldsymbol{\lambda}^*)^\top \mathbf{C}(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \hat{\phi}_t^*(\theta, a)) \right] \\ &\quad + \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \right) \right] + \left(\sum_{\theta \in \Theta^*} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \phi_1^*(\theta, a) \right) \right) \cdot \mathbb{E}[T - T_e] \\ &\quad + (\boldsymbol{\lambda}^*)^\top \mathbb{E} \left[T \boldsymbol{\rho}_1 - \sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} u(\theta) \mathbf{C}(\theta, a) \hat{\phi}_t^*(\theta, a) \right] \\ &\quad + \left(\sum_{\theta \in \Theta, a \in A^+} (u(\theta)(R(\theta, a) - (\boldsymbol{\lambda}^*)^\top \mathbf{C}(\theta, a))) \phi_1^*(\theta, a) \right) \cdot \mathbb{E}[T - T_e] \\ &\stackrel{(e)}{\leq} \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} (u(\theta)(R(\theta, a) - (\boldsymbol{\lambda}^*)^\top \mathbf{C}(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \hat{\phi}_t^*(\theta, a)) \right] \\ &\quad + \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \right) \right] \\ &\quad + (\boldsymbol{\lambda}^*)^\top \mathbb{E}[\mathbf{B}_{T_e+1}] + \max_{\theta \in \Theta, a \in A^+} (R(\theta, a) - (\boldsymbol{\lambda}^*)^\top \mathbf{C}(\theta, a)) \cdot \mathbb{E}[T - T_e]. \end{aligned}$$

565 In the above set of derivations, (a) holds since $T_0 \geq T_e$, (b) is due to Optional Stopping Theorem
 566 since T_e is a stopping time, (c) is by the strong duality of $J(\rho_1)$ as given by (4), (d) establishes by
 567 rearranging terms. At last, for (e), the diminishing term is by strong duality, the transformation from
 568 the fourth term in (d) to the third term in (e) is derived by another application of Optional Stopping
 569 Theorem on the accumulated consumption vector, and for the last term, the upper bound is achieved
 570 since $\sum_{a \in A^+} \phi_1^*(\theta, a) \leq 1$ for any $\theta \in \Theta$ and $\sum_{\theta \in \Theta} u(\theta) = 1$.

571 C.3 Proof of Lemma C.2

572 We will apply the stability result in Chen et al. [2022] as an intermediate to prove our version.
 573 As given, we know that $J(\rho_1)$ and $\hat{J}(\rho_t, \mathcal{H}_t)$ has the same set of basic/non-basic variables and
 574 binding/non-binding constraints as long as the following conditions hold for some constant $D_0 > 0$:

$$\begin{aligned} & \left\| \left(u(\theta) \sum_{\gamma} v(\gamma) r(\theta, a, \gamma) - \hat{u}_t(\theta) \sum_{\gamma} \hat{v}_t(\gamma) r(\theta, a, \gamma) \right)_{(\theta, a) \in \Theta \times A^+} \right\|_{\infty} \leq D_0, \\ & \left\| \left(u(\theta) \sum_{\gamma} v(\gamma) c^i(\theta, a, \gamma) - \hat{u}_t(\theta) \sum_{\gamma} \hat{v}_t(\gamma) c^i(\theta, a, \gamma) \right)_{(\theta, a) \in \Theta \times A^+} \right\|_{\infty} \leq D_0, \quad \forall i \in [n], \\ & \|\rho_1|_S - \rho_t|_S\|_{\infty} \leq D_0, \quad \max \{\rho_1|_{\mathcal{T}} - \rho_t|_{\mathcal{T}}\} \leq D_0. \end{aligned} \quad (11)$$

575 Now, by a standard insertion technique, we have

$$\begin{aligned} & u(\theta) \sum_{\gamma} v(\gamma) r(\theta, a, \gamma) - \hat{u}_t(\theta) \sum_{\gamma} \hat{v}_t(\gamma) r(\theta, a, \gamma) \\ &= (u(\theta) - \hat{u}_t(\theta)) \sum_{\gamma} v(\gamma) r(\theta, a, \gamma) + \hat{u}_t(\theta) \sum_{\gamma} (v(\gamma) - \hat{v}_t(\gamma)) r(\theta, a, \gamma) \\ &\stackrel{(a)}{\leq} \|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_{\infty} + \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1. \end{aligned} \quad (12)$$

576 For (a), the first term is bounded since $r(\theta, a, \gamma) \leq 1$ and $\sum_{\gamma} v(\gamma) = 1$. The second term is similarly
 577 bounded as $\hat{u}_t(\theta) \leq 1$. Therefore, we let $D = D_0/2$, then when we have

$$\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_{\infty} \leq D, \quad \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \leq D,$$

578 the first condition in (11) is met. An almost identical reasoning also holds for the second condition in
 579 (11). Consequently we finish the proof of the lemma.

580 C.4 Proof of Lemma C.3

581 Recall that we are going to prove that

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} (u(\theta) (R(\theta, a) - (\lambda^*)^{\top} C(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \hat{\phi}_t^*(\theta, a)) \right] = O(1), \\ & \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \right) \right] = O(1), \end{aligned}$$

582 when $T \rightarrow \infty$ under Assumption 3.1. For simplicity, we give the following abbreviations:

$$\begin{aligned} P_t &:= \sum_{\theta \in \Theta, a \in A^+} (u(\theta) (R(\theta, a) - (\lambda^*)^{\top} C(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \hat{\phi}_t^*(\theta, a)), \\ Q_t &:= \sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \right), \\ \mathcal{E}_{u,t} &:= [\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_{\infty} \leq D], \quad \mathcal{E}_{v,t} := [\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \leq D]. \end{aligned}$$

On this end, we first utilize Lemma C.2 to show that when condition (6) holds, we have

$$P_t = Q_t = 0.$$

Specifically, for P_t , by the dual feasibility of $J(\rho_1)$, we have $u(\theta) (R(\theta, a) - (\lambda^*)^\top C(\theta, a)) - \mu^*(\theta) \leq 0$. When $u(\theta) (R(\theta, a) - (\lambda^*)^\top C(\theta, a)) - \mu^*(\theta) < 0$, by primal optimality, $\phi_1^*(\theta, a) = 0$ and thus (θ, a) is non-basic for $J(\rho_1)$. By Lemma C.2, (θ, a) is also non-basic for $\hat{J}(\rho_t, \mathcal{H}_t)$ and $\hat{\phi}_t^*(\theta, a) = 0$ holds as well. In conjunction with the case that $u(\theta) (R(\theta, a) - (\lambda^*)^\top C(\theta, a)) - \mu^*(\theta) = 0$, we obtain that $P_t = 0$.

For Q_t , notice that we have $\mu^*(\theta) \geq 0$ for any $\theta \in \Theta$. The case that $\mu^*(\theta) = 0$, again, does not contribute to the total sum. When $\mu^*(\theta) > 0$, by complementary slackness, $\sum_{a \in A^+} \phi_1^*(\theta, a) = 1$, i.e., θ is a binding constraint for $J(\rho_1)$. This, by Lemma C.2, implies that θ is also binding for $\hat{J}(\rho_t, \mathcal{H}_t)$, which shows that the second term is also zero.

With the above, it remains to consider the situation that condition (6) does not hold when $t \leq T_e$, or in other words, $\mathcal{E}_{u,t} \wedge \mathcal{E}_{v,t}$ does not hold. Note that $P_t \leq 1$ and $Q_t \leq 1$ always hold. Thus, we only need to bound the probability that $\neg(\mathcal{E}_{u,t} \wedge \mathcal{E}_{v,t})$. By a union bound, we have

$$\Pr[\neg(\mathcal{E}_{u,t} \wedge \mathcal{E}_{v,t})] = \Pr[\neg\mathcal{E}_{u,t} \vee \neg\mathcal{E}_{v,t}] \leq \Pr[\neg\mathcal{E}_{u,t}] + \Pr[\neg\mathcal{E}_{v,t}].$$

For the first term above, we apply the Hoeffding's inequality and a union bound to derive that

$$\Pr[\neg\mathcal{E}_{u,t}] = \Pr[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_\infty > D] \leq 2|\Theta| \exp(-2D^2(t-1)).$$

Whereas for the second term, we use the concentration result in Weissman et al. [2003] to derive that

$$\Pr[\neg\mathcal{E}_{v,t}] = \Pr[\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 > D] \leq (2^{|\Gamma|} - 2) \exp(-D^2(t-1)/2).$$

Synthesizing the above all, we have

$$\begin{aligned} \mathbb{E}[P_t] &= \mathbb{E}[P_t \mid \mathcal{E}_{u,t} \wedge \mathcal{E}_{v,t}] \cdot \Pr[\mathcal{E}_{u,t} \wedge \mathcal{E}_{v,t}] + \mathbb{E}[P_t \mid \neg(\mathcal{E}_{u,t} \wedge \mathcal{E}_{v,t})] \cdot \Pr[\neg(\mathcal{E}_{u,t} \wedge \mathcal{E}_{v,t})] \\ &\leq 0 + 1 \cdot \Pr[\neg(\mathcal{E}_{u,t} \wedge \mathcal{E}_{v,t})] \\ &\leq 2|\Theta| \exp(-2D^2(t-1)) + (2^{|\Gamma|} - 2) \exp(-D^2(t-1)/2), \end{aligned} \quad (13)$$

$$\mathbb{E}[Q_t] \leq 2|\Theta| \exp(-2D^2(t-1)) + (2^{|\Gamma|} - 2) \exp(-D^2(t-1)/2). \quad (14)$$

Summing (13) and (14) from 1 to T_e , we achieve that

$$\begin{aligned} &\left\{ \mathbb{E} \left[\sum_{t=1}^{T_e} P_t \right], \mathbb{E} \left[\sum_{t=1}^{T_e} Q_t \right] \right\} \\ &\leq \sum_{t=1}^{T_e} \left(2|\Theta| \exp(-2D^2(t-1)) + (2^{|\Gamma|} - 2) \exp(-D^2(t-1)/2) \right) \\ &\leq \frac{2|\Theta|}{1 - \exp(-2D^2)} + \frac{2^{|\Gamma|} - 2}{1 - \exp(-D^2/2)}, \end{aligned}$$

which conclude the proof of the lemma.

C.5 Proof of Lemma C.4

As implied by (10), the proof of this lemma reduces to bound $\mathbb{E}[T - T_e]$, i.e., showing that T_e is sufficiently close to T . On this side, we first recall the definition of T_e in (9):

$$T_e := \min\{T_0, \min\{t : \max\{\|\rho_1|_S - \rho_t|_S\|_\infty, \max\{\rho_1|_{\mathcal{T}} - \rho_t|_{\mathcal{T}}\}\} > D\} - 1\},$$

where T_0 is the stopping time of Algorithm 1, and S and \mathcal{T} correspondingly represent the set of binding/non-binding resource constraints in LP $J(\rho_1)$. For simplicity, we define

$$\mathcal{N}(\rho_1, D, S) := \{\kappa : \max\{\|\rho_1|_S - \kappa|_S\|_\infty, \max\{\rho_1|_{\mathcal{T}} - \kappa|_{\mathcal{T}}\}\} \leq D\}.$$

606 It is without loss of generality to suppose that $D < \rho^{\min}$. We let

$$T_D := \min\{t : \rho_t \notin \mathcal{N}(\rho_1, D, \mathcal{S})\} - 1, \quad T_- = \lfloor T + 1 - 1/(\rho^{\min} - D) \rfloor.$$

607 We show that if $t \leq T_-$ and $t \leq T_D$, then $t \leq T_e$. In fact, under the condition, we derive that

$$B_t \geq (T - t + 1)(\rho_1 - D\mathbf{1}) \geq \frac{1}{\rho^{\min} - D}(\rho_1 - D\mathbf{1}) \geq \mathbf{1},$$

608 which implies that $t \leq T_0$, and therefore $t \leq T_e$. As a result, we have

$$\mathbb{E}[T_e] = \sum_{t=1}^T \Pr[T_e \geq t] \geq \sum_{t=1}^{T_-} \Pr[T_e \geq t] \geq \sum_{t=1}^{T_-} \Pr[T_D \geq t] = T_- - \sum_{t=1}^{T_-} \Pr[t > T_D]. \quad (15)$$

609 Before we continue to bound (15), we first give an observation on the dynamics of ρ_t . By the update
610 process of the budget, we have for any $t \geq 1$,

$$\begin{aligned} B_{t+1} = B_t - c_t &\implies \rho_{t+1}(T - t) = \rho_t(T - t + 1) - c_t \\ &\implies \rho_{t+1} = \rho_t + \frac{\rho_t - c_t}{T - t}. \end{aligned}$$

611 Now let

$$M_t^C := \frac{\rho_t - \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) C(\theta, a) \right]}{T - t}, \quad N_t^C := \frac{\mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) C(\theta, a) \right] - c_t}{T - t}.$$

612 We then have

$$\rho_{t+1} - \rho_t = \frac{\rho_t - c_t}{T - t} = M_t^C + N_t^C. \quad (16)$$

613 We now define an auxiliary process which benefits the analysis. Specifically, for $t \in [T]$, let

$$\tilde{\rho}_t := \begin{cases} \rho_t, & t \leq T_D; \\ \rho_{T_D}, & t > T_D. \end{cases}$$

614 Therefore,

$$\tilde{\rho}_{t+1} - \tilde{\rho}_t = \begin{cases} M_t^C + N_t^C, & t \leq T_D; \\ 0, & t > T_D. \end{cases}$$

615 We further define the following two auxiliary variables for $t \in [T]$:

$$\tilde{M}_t^C := \begin{cases} M_t^C, & t \leq T_D; \\ 0, & t > T_D. \end{cases}, \quad \tilde{N}_t^C := \begin{cases} N_t^C, & t \leq T_D; \\ 0, & t > T_D. \end{cases}$$

616 As a result, we have

$$\tilde{\rho}_{t+1} - \tilde{\rho}_t = \tilde{M}_t^C + \tilde{N}_t^C.$$

617 Now we come back to (15). Notice that

$$\begin{aligned} &\Pr[t > T_D] \\ &= \Pr[\rho_s \notin \mathcal{N}(\rho_1, D, \mathcal{S}) \text{ for some } s \leq t] = \Pr[\tilde{\rho}_t \notin \mathcal{N}(\rho_1, D, \mathcal{S})] \\ &\leq \Pr \left[\left\| \sum_{\tau=1}^{t-1} (\tilde{M}_\tau^C + \tilde{N}_\tau^C) \right\|_{\mathcal{S}} \geq D \text{ or } \min_{\tau=1}^{t-1} \left(\tilde{M}_\tau^C + \tilde{N}_\tau^C \right) |_{\mathcal{T}} < -D \right] \\ &\leq \Pr \left[\left\| \sum_{\tau=1}^{t-1} \tilde{M}_\tau^C \right\|_{\mathcal{S}} \geq D/2 \text{ or } \min_{\tau=1}^{t-1} \tilde{M}_\tau^C |_{\mathcal{T}} < -D/2 \right] + \Pr \left[\left\| \sum_{\tau=1}^{t-1} \tilde{N}_\tau^C \right\|_{\infty} \geq D/2 \right]. \end{aligned} \quad (18)$$

618 For the second term in (18), we observe that each entry of $\{\sum_{\tau < t} \tilde{N}_\tau^C\}_t$ is a martingale with the
619 absolute value of the τ -th increment bounded by $1/(T - \tau)$. Since

$$\sum_{\tau=1}^{t-1} \frac{1}{(T - \tau)^2} \leq \frac{1}{T - t},$$

by applying the Azuma–Hoeffding inequality and a union bound, we achieve that

$$\Pr \left[\left\| \sum_{\tau=1}^{t-1} \widetilde{\mathbf{M}}_{\tau}^C \right\|_{\infty} \geq D/2 \right] \leq 2n \exp \left(-\frac{(T-t)D^2}{8} \right).$$

We now come back to the first term in (18), for any $\{D_1, \dots, D_{t-1}\}$ such that $\sum_{\tau=1}^{t-1} D_{\tau}/(T-\tau) \leq D/2$, we have

$$\begin{aligned} & \left\{ \left\| \sum_{\tau=1}^{t-1} \widetilde{\mathbf{M}}_{\tau}^C \right\|_{\infty} > D/2 \text{ or } \min \sum_{\tau=1}^{t-1} \widetilde{\mathbf{M}}_{\tau}^C|_{\mathcal{T}} < -D/2 \right\} \\ \Rightarrow & \left\{ \left\| \widetilde{\mathbf{M}}_{\tau}^C \right\|_{\infty} > \frac{D_{\tau}}{T-\tau} \text{ or } \min \widetilde{\mathbf{M}}_{\tau}^C|_{\mathcal{T}} < -\frac{D_{\tau}}{T-\tau} \right\} \text{ for some } \tau \in [T-1]. \end{aligned}$$

We now define

$$\mathcal{E}_{\tau}(D_{\tau}) := \left(\left\| \mathbf{M}_{\tau}^C \right\|_{\mathcal{S}} \leq \frac{D_{\tau}}{T-\tau} \right) \wedge \left(\min \mathbf{M}_{\tau}^C|_{\mathcal{T}} \geq -\frac{D_{\tau}}{T-\tau} \right) \text{ holds for } \forall \rho_{\tau} \in \mathcal{N}(\rho_1, D, \mathcal{S}).$$

Since $\widetilde{\mathbf{M}}_{\tau}^C \neq 0$ only when $t \leq T_D$, i.e., $\rho_t \in \mathcal{N}(\rho_1, D, \mathcal{S})$, by the definition of $\mathcal{E}_{\tau}(D_{\tau})$, we have the following claim:

$$\left\{ \left\| \sum_{\tau=1}^{t-1} \widetilde{\mathbf{M}}_{\tau}^C \right\|_{\infty} > D/2 \text{ or } \min \sum_{\tau=1}^{t-1} \widetilde{\mathbf{M}}_{\tau}^C|_{\mathcal{T}} < -D/2 \right\} \subseteq \bigcup_{\tau=1}^{t-1} \neg \mathcal{E}_{\tau}(D_{\tau}), \quad \forall \sum_{\tau=1}^{t-1} \frac{D_{\tau}}{T-\tau} \leq D/2. \quad (19)$$

Thus, we forward to bound $\Pr[\neg \mathcal{E}_{\tau}(D_{\tau})]$ for a suitable choice of $\{D_{\tau}\}_{1 \leq \tau \leq T}$. Recall that we have defined events $\mathcal{E}_{u,\tau}$ and $\mathcal{E}_{v,\tau}$ as follows:

$$\mathcal{E}_{u,\tau} := [\|(u(\theta) - \widehat{u}_{\tau}(\theta))_{\theta \in \Theta}\|_{\infty} \leq D], \quad \mathcal{E}_{v,\tau} := [\|(v(\gamma) - \widehat{v}_{\tau}(\gamma))_{\gamma \in \Gamma}\|_1 \leq D].$$

We have the following lemma, which we are going to prove in Appendix C.6:

Lemma C.5. When $\rho_{\tau} \in \mathcal{N}(\rho_1, D, \mathcal{S})$ and $\mathcal{E}_{u,\tau} \wedge \mathcal{E}_{v,\tau}$ hold,

$$\begin{aligned} (T-\tau) \left\| \mathbf{M}_{\tau}^C \right\|_{\mathcal{S}} & \leq \|(u(\theta) - \widehat{u}_{\tau}(\theta))_{\theta \in \Theta}\|_1 + \|(v(\gamma) - \widehat{v}_{\tau}(\gamma))_{\gamma \in \Gamma}\|_1, \\ (T-\tau) \min \mathbf{M}_{\tau}^C|_{\mathcal{T}} & \geq -\|(u(\theta) - \widehat{u}_{\tau}(\theta))_{\theta \in \Theta}\|_1 - \|(v(\gamma) - \widehat{v}_{\tau}(\gamma))_{\gamma \in \Gamma}\|_1. \end{aligned}$$

Further, it is clear that $(T-\tau) \left\| \mathbf{M}_{\tau}^C \right\|_{\mathcal{S}} \leq 1$ and $(T-\tau) \min \mathbf{M}_{\tau}^C|_{\mathcal{T}} \geq -1$ holds. Inspired by the above observations, we let the series of D_1, \dots, D_{T-1} be the following form:

$$D_{\tau} = \begin{cases} 1, & \tau \leq \eta T; \\ (\tau-1)^{-1/4}, & \tau > \eta T, \end{cases}$$

where $\eta \in (0, 1)$ is a constant to be specified. We need to satisfy the following constraints:

$$\sum_{t=1}^{T-1} \frac{D_t}{T-t} \leq D/2, \quad (\eta T)^{-1/4} < D.$$

Here, the first constraint is instructed by (19), and the second is to guarantee that when $\|(u(\theta) - \widehat{u}_{\tau}(\theta))_{\theta \in \Theta}\|_1 + \|(v(\gamma) - \widehat{v}_{\tau}(\gamma))_{\gamma \in \Gamma}\|_1 < (\tau-1)^{-1/4}$ for $\tau > \eta T$, $\mathcal{E}_{u,\tau} \wedge \mathcal{E}_{v,\tau}$ naturally holds, and therefore we can apply Lemma C.5. For the first one, we notice that

$$\sum_{\tau=1}^{T-1} \frac{D_{\tau}}{T-\tau} = \sum_{\tau=1}^{\eta T} \frac{1}{T-\tau} + \sum_{\tau=\eta T+1}^{T-1} \frac{1}{(T-\tau)(\tau-1)^{1/4}} \leq \log \frac{T-1}{(1-\eta)T-1} + \frac{\log T}{(\eta T)^{1/4}}.$$

Therefore, for some η such that $\log(1 - \eta) \geq -D/4$, $\sum_{t=1}^T D_t/(T - t) \leq D/2$ establishes for sufficiently large $T \gg 1$, and the second constraint is also satisfied.

We are now prepared to bound $\Pr[\neg \mathcal{E}_\tau(D_\tau)]$ for the $\{D_\tau\}$ we just proposed. To start with, when $\tau \leq \eta T$, $\mathcal{E}_\tau(D_\tau)$ always holds, thus $\Pr[\neg \mathcal{E}_\tau(D_\tau)] = 0$. When $\tau > \eta T$, since $\tau^{-1/4}/2 < D$, by Hoeffding's inequality and union bound, we have

$$\begin{aligned} & \Pr[\neg \mathcal{E}_\tau(D_\tau)] \\ & \stackrel{(a)}{\leq} \Pr \left[\|(u(\theta) - \hat{u}_\tau(\theta))_{\theta \in \Theta}\|_1 \leq (\tau - 1)^{-1/4}/2 \right] + \Pr \left[\|(v(\gamma) - \hat{v}_\tau(\gamma))_{\gamma \in \Gamma}\|_1 \leq (\tau - 1)^{-1/4}/2 \right] \\ & \leq 2|\Theta| \exp \left(-\frac{(\tau - 1)^{1/2}}{8|\Theta|^2} \right) + 2|\Gamma| \exp \left(-\frac{(\tau - 1)^{1/2}}{8|\Gamma|^2} \right). \end{aligned}$$

Here, (a) is by Lemma C.5 and a union bound. Therefore, according to (19), we have

$$\Pr \left[\left\| \sum_{\tau=1}^{t-1} \tilde{M}_\tau^C \right\|_S > D/2 \text{ or } \min_{\tau=1}^{t-1} \tilde{M}_\tau^C|_{\mathcal{T}} < -D/2 \right] \leq \sum_{\tau=1}^{t-1} \Pr[\neg \mathcal{E}_\tau(D_\tau)],$$

and therefore,

$$\Pr \left[\left\| \sum_{\tau=1}^{t-1} \tilde{M}_\tau^C \right\|_S > D/2 \text{ or } \min_{\tau=1}^{t-1} \tilde{M}_\tau^C|_{\mathcal{T}} < -D/2 \right] \leq \begin{cases} 0, & t \leq \eta T + 1; \\ \sum_{\tau=\eta T+1}^{t-1} \exp \left\{ -\tau^{1/2} \right\}, & t > \eta T + 1. \end{cases}$$

Plugging the into (18) and (15), we obtain that when $T \rightarrow \infty$,

$$\begin{aligned} & \mathbb{E}[T - T_e] \\ & \leq T - T_- \\ & \quad + \sum_{t=1}^{T_-} \left(\Pr \left[\left\| \sum_{\tau=1}^{t-1} \tilde{M}_\tau^C \right\|_S > D/2 \text{ or } \min_{\tau=1}^{t-1} \tilde{M}_\tau^C|_{\mathcal{T}} < -D/2 \right] + 2n \exp \left(-\frac{(T - t)D^2}{8} \right) \right) \\ & \leq \frac{1}{\rho_{\min} - D} + 2n(1 - \exp(-D^2/8))^{-1} + O(T^2) \exp(-T^{1/2}) = O(1). \end{aligned}$$

At last, combining with (10), we finally finish the proof of Lemma C.4.

C.6 Proof of Lemma C.5

To start with, we notice that

$$(T - \tau)M_\tau^C = \rho_\tau - \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) C(\theta, a) \right].$$

Now, notice that $\rho_\tau \in \mathcal{N}(\rho_1, D, \mathcal{S})$ and $\mathcal{E}_{u,\tau} \wedge \mathcal{E}_{v,\tau}$ are the condition of Lemma C.2, therefore, the set of resource binding constraints of $\hat{J}(\rho_t, \mathcal{H}_t)$ are identical to that of $J(\rho_1)$, i.e., \mathcal{S} . Hence, for any $i \in [n]$,

$$\begin{aligned} & \rho_\tau^i|_{\mathcal{S}} - \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) C^i(\theta, a)|_{\mathcal{S}} \right] \\ & = \sum_{\theta \in \Theta, a \in A^+} \hat{u}_\tau(\theta) \hat{\phi}_\tau^*(\theta, a) \sum_{\gamma} \hat{v}_\tau(\gamma) c^i(\theta, a, \gamma)|_{\mathcal{S}} - \sum_{\theta \in \Theta, a \in A^+} u(\theta) \hat{\phi}_\tau^*(\theta, a) \sum_{\gamma} v(\gamma) c^i(\theta, a, \gamma)|_{\mathcal{S}} \\ & = \sum_{\theta \in \Theta, a \in A^+} (u(\theta) - \hat{u}_\tau(\theta)) \hat{\phi}_\tau^*(\theta, a) \sum_{\gamma} v(\gamma) c^i(\theta, a, \gamma)|_{\mathcal{S}} \\ & \quad + \sum_{\theta \in \Theta, a \in A^+} \hat{u}_\tau(\theta) \hat{\phi}_\tau^*(\theta, a) \sum_{\gamma} (\hat{v}_\tau(\gamma) - v(\gamma)) c^i(\theta, a, \gamma)|_{\mathcal{S}} \\ & \stackrel{(a)}{\leq} \|(u(\theta) - \hat{u}_\tau(\theta))_{\theta \in \Theta}\|_1 + \|(v(\gamma) - \hat{v}_\tau(\gamma))_{\gamma \in \Gamma}\|_1. \end{aligned}$$

Here, the bound on the first term in (a) establishes because for any $\theta \in \Theta$,

$$\sum_{a \in A^+} \widehat{\phi}_\tau^*(\theta, a) \sum_{\gamma} v(\gamma) \mathbf{c}^i(\theta, a, \gamma) |_{\mathcal{S}} \leq 1$$

since $\sum_{a \in A^+} \widehat{\phi}_\tau^*(\theta, a) \leq 1$. The bound on the second term is similar. Thus, we achieve the result for binding constraints. The proof for non-binding constraints resembles the above by noticing that

$$\boldsymbol{\rho}_\tau |_{\mathcal{T}} \geq \sum_{\theta \in \Theta, a \in A^+} \widehat{u}_\tau(\theta) \widehat{\phi}_\tau^*(\theta, a) \sum_{\gamma} \widehat{v}_\tau(\gamma) \mathbf{c}(\theta, a, \gamma) |_{\mathcal{T}}.$$

D Missing Proofs in Section 4

D.1 Proof of Theorem 4.1

With Lemma 4.1 in hand, we now show how to derive Theorem 4.1. Specifically, the regret decomposition technique in Lemma C.1 still works fine. We only need to re-derive corresponding results for Lemmas C.3 and C.4. We have the following results on this side, which are proved respectively in Appendices D.3 and D.4.

Lemma D.1. *Under Assumption 3.1, with partial information feedback, we have when $T \rightarrow \infty$:*

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta, a \in A^+} (u(\theta) (R(\theta, a) - (\boldsymbol{\lambda}^*)^\top \mathbf{C}(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \widehat{\phi}_t^*(\theta, a)) \right] &= O(1), \\ \mathbb{E} \left[\sum_{t=1}^{T_e} \sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \widehat{\phi}_t^*(\theta, a) \right) \right] &= O(\log T). \end{aligned}$$

Lemma D.2. *Under Assumption 3.1, with partial information feedback, we have when $T \rightarrow \infty$:*

$$(\boldsymbol{\lambda}^*)^\top \mathbb{E} [\mathbf{B}_{T_e+1}] + \max_{\theta \in \Theta, a \in A^+} (R(\theta, a) - (\boldsymbol{\lambda}^*)^\top \cdot \mathbf{C}(\theta, a)) \mathbb{E} [T - T_e] = O(1).$$

Lemmas C.1, D.1 and D.2 in together leads to Theorem 4.1.

D.2 Proof of Lemma 4.1

Some preparations are required before we come to prove the lemma. To start with, we notice that $Y_\tau = \Pr[a_1 \neq 0] + \dots + \Pr[a_{t-1} \neq 0]$. By the control rule of Algorithm 1, we have

$$\Pr[a_\tau \neq 0] = \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \widehat{\phi}_\tau^*(\theta, a) \mid \mathcal{H}_\tau \right].$$

We first give a lower bound on $\mathbb{E}_{\theta \sim \mathcal{U}} [\sum_{a \in A^+} \widehat{\phi}_\tau^*(\theta, a) \mid \mathcal{H}_\tau]$ with $\boldsymbol{\rho}_\tau$, taking $\mathbb{E}_{\theta \sim \widehat{\mathcal{U}}_\tau} [\sum_{a \in A^+} \widehat{\phi}_\tau^*(\theta, a) \mid \mathcal{H}_\tau]$ as an intermediate.

Lemma D.3.

$$\mathbb{E}_{\theta \sim \widehat{\mathcal{U}}_\tau} \left[\sum_{a \in A^+} \widehat{\phi}_\tau^*(\theta, a) \mid \mathcal{H}_\tau \right] \geq \min \{1, \min \boldsymbol{\rho}_\tau\}.$$

Proof of Lemma D.3. To start with, when $\boldsymbol{\rho}_\tau \geq \mathbf{1}$, then clearly, all the resource constraints in $\widehat{J}(\boldsymbol{\rho}_\tau, \mathcal{H}_\tau)$ are satisfied even when $\sum_{a \in A^+} \phi(\theta, a) = 1$ holds for any $\theta \in \Theta$. Therefore, an optimal solution should have this form.

667 We now consider the case that $\min \rho_\tau < 1$. In this case, if there is a feasible solution that
 668 $\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) = 1$ holds for any $\theta \in \Theta$, then the proof is also finished. Otherwise, there is
 669 at least a binding resource constraint in $\hat{J}(\rho_\tau, \mathcal{H}_\tau)$, which we denote by i^* . Consequently,

$$\mathbb{E}_{\theta \sim \hat{\mathcal{U}}_\tau} \left[\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \right] \geq \mathbb{E}_{\theta \sim \hat{\mathcal{U}}_\tau} \left[\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \hat{C}_\tau^{i^*}(\theta, a) \right] = \rho_\tau^{i^*} \geq \min \rho_\tau.$$

670 This finishes the proof of the lemma. \square

671 Thus, we have

$$\begin{aligned} \Pr[a_\tau \neq 0] &= \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \mid \mathcal{H}_\tau \right] \geq \mathbb{E}_{\theta \sim \hat{\mathcal{U}}_\tau} \left[\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \mid \mathcal{H}_\tau \right] - \|u(\theta) - \hat{u}_\tau(\theta)\|_1 \\ &\geq \min \{1, \min \rho_\tau\} - \|u(\theta) - \hat{u}_\tau(\theta)\|_1. \end{aligned} \quad (20)$$

672 Further, we have the following result bounding $\min \rho_\tau$ when t is no larger than a fraction of T .

673 **Lemma D.4.** When $t \leq (\rho^{\min}/2) \cdot T$, $\min \rho_\tau \geq \rho^{\min}/2$.

674 *Proof of Lemma D.4.* In fact, for $t \leq (\rho^{\min}/2) \cdot T$,

$$\rho_\tau = \frac{T \cdot \rho_1 - \sum_{\tau=1}^{t-1} c_\tau}{T - t + 1} \geq \frac{T \cdot \rho_1 - t \cdot 1}{T} \geq \frac{\rho_1}{2}.$$

675 This concludes the proof. \square

676 Now, by Weissman et al. [2003], with probability $1 - O(1/T)$, we have

$$\|u(\theta) - \hat{u}_\tau(\theta)\|_1 \leq \frac{\rho^{\min}}{4}, \quad \forall \tau \geq \Theta(\log T).$$

677 Taking into (20), we derive that

$$\Pr[a_\tau \neq 0] \geq \frac{\rho^{\min}}{4}, \quad \forall \Theta(\log T) \leq t \leq \frac{\rho^{\min}}{2} \cdot T.$$

678 Consequently, within the period, the probability that there are $\Omega(\log T)$ consecutive rounds in which
 679 the agent chooses to quit in all these rounds is $O(1/T)$. This proves the first part. Meanwhile, at time
 680 $t = \lceil (\rho^{\min}/2) \cdot T \rceil + 1$, by Azuma–Hoeffding inequality, we derive that with probability $1 - O(1/T)$,
 681 $Y_t = \sum_{\tau=1}^{t-1} \Pr[a_\tau \neq 0] \geq \Omega(T)$, which proves the second part.

682 D.3 Proof of Lemma D.1

683 We concentrate on adapting the proof of Lemma C.3 into the partial information feedback setting.
 684 To start with, we suppose that the conditions given in Lemma 4.1 hold. In fact, since the failure
 685 probability is only $O(1/T)$, and the sum is upper bounded by $O(T)$, therefore the failure case only
 686 contributes $O(1)$ to the total expectation.

687 Now, recall the following definitions:

$$P_t := \sum_{\theta \in \Theta, a \in A^+} (u(\theta) (R(\theta, a) - (\lambda^*)^\top C(\theta, a)) - \mu^*(\theta)) (\phi_1^*(\theta, a) - \hat{\phi}_t^*(\theta, a)),$$

$$Q_t := \sum_{\theta \in \Theta} \mu^*(\theta) \left(1 - \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \right),$$

$$\mathcal{E}_{u,t} := [\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_\infty \leq D], \quad \mathcal{E}_{v,t} := [\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \leq D],$$

688 and by (13) and (14), we have

$$\mathbb{E}[P_t] \leq \Pr[\neg \mathcal{E}_{u,t}] + \Pr[\neg \mathcal{E}_{v,t}], \quad \mathbb{E}[Q_t] \leq \Pr[\neg \mathcal{E}_{u,t}] + \Pr[\neg \mathcal{E}_{v,t}].$$

Now, the bound on $\Pr[\neg \mathcal{E}_{u,t}]$ inherits the analysis in the proof of Lemma C.3, as partial information feedback does not affect the learning of the request distribution. That is,

$$\Pr[\neg \mathcal{E}_{u,t}] = \Pr[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_\infty > D] \leq 2|\Theta| \exp(-2D^2(t-1)).$$

For $\Pr[\neg \mathcal{E}_{v,t}]$, when $t \leq \Theta(\log T)$, it is obviously bounded by 1. By Lemma 4.1, when $\Theta(\log T) \leq t \leq C_b \cdot T$, by Weissman et al. [2003], we have

$$\Pr[\neg \mathcal{E}_{v,t}] = \Pr[\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 > D] \leq (2^{|\Gamma|} - 2) \exp\left(-\frac{D^2 C_f(t-1)}{2 \log T}\right).$$

Further, when $t > C_b \cdot T$, we correspondingly derive

$$\Pr[\neg \mathcal{E}_{v,t}] \leq (2^{|\Gamma|} - 2) \exp\left(-\frac{D^2 C_r(t-1)}{2}\right).$$

Putting the above together, we achieve that

$$\begin{aligned} & \left\{ \mathbb{E} \left[\sum_{t=1}^{T_e} P_t \right], \mathbb{E} \left[\sum_{t=1}^{T_e} Q_t \right] \right\} \\ & \leq \sum_{t=1}^T 2|\Theta| \exp(-2D^2(t-1)) + \Theta(\log T) \\ & \quad + (2^{|\Gamma|} - 2) \left(\sum_{t=\Theta(\log T)}^{C_b \cdot T} \exp\left(-\frac{D^2 C_f(t-1)}{2 \log T}\right) + \sum_{t=C_b \cdot T+1}^T \exp\left(-\frac{D^2 C_r(t-1)}{2}\right) \right) \\ & \leq \Theta(1) + \Theta(\log T) + (2^{|\Gamma|} - 2) \left(\frac{\Theta(1)}{1 - \exp(-\Theta(1/\log T))} + \exp(-\Theta(T)) \right) \\ & \leq O(\log T). \end{aligned}$$

Here, for the last inequality, by Taylor expansion, we have $1 - e^{-x} \geq x - x^2/2$ for $x > 0$, therefore,

$$\frac{1}{1 - \exp(-\Theta(1/\log T))} \leq \frac{1}{\Theta(1/\log T) - 1/\log^2 T} \leq \Theta(\log T).$$

This finishes the proof.

D.4 Proof of Lemma D.2

As in Appendix D.3 when we prove Lemma C.3, we only consider the case when the conditions in Lemma 4.1 establish, as the contribution of the failure cases on the expectation-sum is $O(1)$. We now bound $\mathbb{E}[T - T_e]$ in the good case when the sample accessing frequency under partial information feedback is guaranteed. Specifically, as predefined in the proof of Lemma C.3, we only need to re-calculate the following, as the other terms remain unchanged with partial information:

$$\sum_{t=\eta T+2}^{T_-} \sum_{\tau=\eta T+1}^{t-1} \Pr \left[\|(v(\gamma) - \hat{v}_\tau(\gamma))_{\gamma \in \Gamma}\|_1 \leq (\tau-1)^{-1/4}/2 \right].$$

Here, η is specified in the definition of D_τ . It is hard for us to directly compare η and C_b in Lemma 4.1. Nevertheless, in any case, we know that when T is sufficiently large, $Y_\tau/(\tau-1) = \Omega(1/\log T)$ for $\tau \geq \eta T$. Therefore, we have

$$\Pr \left[\|(v(\gamma) - \hat{v}_\tau(\gamma))_{\gamma \in \Gamma}\|_1 \leq (\tau-1)^{-1/4}/2 \right] \leq 2|\Gamma| \exp\left(-\frac{(\tau-1)^{1/2}}{|\Gamma|^2 O(\log T)}\right).$$

706 Hence,

$$\begin{aligned}
& \sum_{t=\eta T+2}^{T_-} \sum_{\tau=\eta T+1}^{t-1} \Pr \left[\|(v(\gamma) - \hat{v}_\tau(\gamma))_{\gamma \in \Gamma}\|_1 \leq (\tau-1)^{-1/4}/2 \right] \\
& \leq 2|\Gamma| \sum_{t=\eta T+2}^{T_-} \sum_{\tau=\eta T+1}^{t-1} \exp \left(-\frac{(\tau-1)^{1/2}}{|\Gamma|^2 O(\log T)} \right) \\
& = O(T^2) \exp \left(-\Omega \left(\frac{T^{1/2}}{\log T} \right) \right) = O(1).
\end{aligned}$$

707 Combining with the other parts, Lemma D.2 is proved.

708 E Missing Proofs in Section 5

709 E.1 Proof of Theorem 5.1

710 We will prove Theorem 5.1 in the following, and we are inspired by the analysis in Chen et al. [2022].

711 E.1.1 Another Regret Decomposition

712 Different from our analysis for the regular cases, in general circumstances, we introduce another regret
713 decomposition method. The reason for involving such an alternative is that without the regularity
714 assumptions, we no longer have any local stability guarantee even when the estimates are close.
715 Therefore, the decision given by Algorithm 1 does not coincides with the optimal decision even when
716 the distribution learning process converges well, and the corresponding analysis in Section 3 does not
717 work out anymore.

718 We now present a more general regret decomposition as follows:

$$\begin{aligned}
V^{\text{FL}} - \text{Rew} &= T \cdot J(\boldsymbol{\rho}_1) - \mathbb{E} \left[\sum_{t=1}^{T_0} r(\theta_t, a_t, \gamma_t) \right] \\
&\stackrel{(a)}{=} T \cdot J(\boldsymbol{\rho}_1) - \mathbb{E} \left[\sum_{t=1}^{T_0} \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \right] \\
&\stackrel{(b)}{=} J(\boldsymbol{\rho}_1) \cdot \mathbb{E} [T - T_0] + \mathbb{E} \left[\sum_{t=1}^{T_0} \left(J(\boldsymbol{\rho}_1) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \right) \right].
\end{aligned} \tag{21}$$

719 Here, (a) holds due to the Optimal Stopping Theorem, since T_0 is a stopping time. Meanwhile, by
720 the decision process, we have for any θ_t :

$$\mathbb{E}_{a_t, \gamma_t} [r(\theta_t, a_t, \gamma_t) \mid \theta_t] = \sum_{a \in A^+} \hat{\phi}_t^*(\theta_t, a) R(\theta_t, a).$$

721 Further, (b) is by a re-arrangement. To give a bound for (21), we respectively analyze $\mathbb{E}[T - T_0]$ the
722 stopping time, and difference between the optimal accumulated rewards and the real ones.

723 E.1.2 Bounding the Stopping Time

724 To settle the stopping time, we first reduce it to $\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)$ for $t \leq T_0$, and then deals with
725 these values. We notice that $t \leq T_0$ as long as that $\mathbf{B}_t \geq \mathbf{1}$, or $\boldsymbol{\rho}_t \geq \mathbf{1}/(T - t + 1)$. Now, since for
726 any $i \in [n]$,

$$\rho_t^i = \rho_1^i - (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t)^i \geq \rho^{\min} - \max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0),$$

727 we have $\min \boldsymbol{\rho}_t \geq \rho^{\min} - \max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)$. Therefore,

$$\begin{aligned}
t \leq T_0 &\iff \rho^{\min} - \max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0) \geq \frac{1}{T - t + 1} \\
&\iff \max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0) \leq \rho^{\min} - \frac{1}{T - t + 1}.
\end{aligned} \tag{22}$$

728 Since $\mathbb{E}[T_0] \geq \Pr[T_0 \geq t] \cdot t$ for any $t \in [T]$, we only need to bound the following term for some
 729 certain t :

$$\Pr \left[\max(\rho_1 - \rho_t, 0) \leq \rho^{\min} - \frac{1}{T-t+1} \right].$$

730 We will further prove the following lemma in Appendix E.3:

731 **Lemma E.1.** *It holds for any $t < T$ that*

$$\Pr \left[\max(\rho_1 - \rho_t, 0) \geq \Theta \left(\frac{1}{T-1} + \sum_{\tau=2}^{t-1} \sqrt{\frac{\log T}{(T-\tau)^2(\tau-1)}} + \sqrt{\frac{\log T}{T-t}} \right) \right] \leq O\left(\frac{1}{T}\right).$$

With the light of Lemma E.1, it is natural for us to compute

$$\sum_{\tau=2}^{t-1} \sqrt{\frac{\log T}{(T-\tau)^2(\tau-1)}} \leq \begin{cases} \sqrt{\log T} \cdot \frac{4\sqrt{t-2}}{T-1}, & 2 \leq t \leq (T+1)/2; \\ \sqrt{\log T} \cdot \frac{2}{\sqrt{T-t}}, & t > (T+1)/2. \end{cases}$$

732 In fact, to derive the above, we notice that when $2 \leq t \leq (T+1)/2$,

$$\sum_{\tau=1}^{t-1} \frac{1}{(T-\tau)(\tau-1)^{1/2}} \leq \frac{2}{T-1} \sum_{\tau=2}^{t-1} \frac{1}{(\tau-1)^{1/2}} \leq \frac{4\sqrt{t-1}}{T-1}.$$

733 Meanwhile, when $t > (T+1)/2$, we have $T-t < t-1$, which leads to

$$\sum_{\tau=2}^{t-1} \frac{1}{(T-\tau)(\tau-1)^{1/2}} \leq \sqrt{\frac{8}{T-1}} + \sum_{\tau=(T+1)/2}^{t-1} \frac{1}{(T-\tau)^{3/2}} \leq \frac{2}{\sqrt{T-t}}.$$

734 With these calculations, we come back to the bound on $\mathbb{E}[T_0]$, we notice that when T is sufficiently
 735 large and $t = T - O(\log T)$, it holds that

$$\Theta \left(\frac{1}{T-1} + \sum_{\tau=2}^{t-1} \sqrt{\frac{\log T}{(T-\tau)^2(\tau-1)}} + \sqrt{\frac{\log T}{T-t}} \right) + \frac{1}{T-t+1} = O(1) \leq \rho^{\min}.$$

736 Thus, we have

$$\begin{aligned} \mathbb{E}[T - T_0] &= T - \mathbb{E}[T_0] \stackrel{(a)}{\leq} T - \Pr[T_0 \geq T - O(\log T)] \cdot (T - O(\log T)) \\ &\stackrel{(b)}{\leq} T - \left(1 - O\left(\frac{1}{T}\right)\right) \cdot (T - O(\log T)) = O(\log T). \end{aligned} \quad (23)$$

737 In the above, (a) is because $\mathbb{E}[T_0] \geq \Pr[T_0 \geq t] \cdot t$ for any fixed t , and (b) is due to Lemma E.1.
 738 Consequently, we finish the analysis of the stopping time in (21).

739 E.1.3 The Gap to the Optimal Reward

740 The rest part of (21) that we are left to consider is the following:

$$\begin{aligned} J(\rho_1) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \\ = \left(J(\rho_1) - \hat{J}(\rho_t, \mathcal{H}_t) \right) + \left(\hat{J}(\rho_t, \mathcal{H}_t) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \right). \end{aligned} \quad (24)$$

741 Note that the second difference term in (24) reflects the estimation error on distributions of the context
 742 and the external factor, which leads to the following result as to be proved in Appendix E.4:

743 **Lemma E.2.** *We have for $t \geq 2$:*

$$\mathbb{E} \left[\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \right] \leq O \left(\sqrt{\frac{\log T}{t-1}} + \frac{1}{T} \right).$$

744 Lemma E.2 induces an $O(\sqrt{T \log T})$ accumulated regret considering (24) when summing from $t = 2$
 745 to $T_0 \leq T$. While for the first term in (24), our main thread here is to bound $\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)$ with $J(\boldsymbol{\rho}_t)$.
 746 To fix the idea, we compare these two optimization problems:

$$\begin{aligned} J(\boldsymbol{\rho}_t) &:= \max_{\phi: \Theta \times A^+ \rightarrow \mathbb{R}_+} \sum_{\theta \in \Theta, a \in A^+} u(\theta) \phi(\theta, a) \sum_{\gamma} r(\theta, a, \gamma) v(\gamma), \\ \text{s.t.} \quad &\sum_{\theta \in \Theta, a \in A^+} u(\theta) \phi(\theta, a) \sum_{\gamma} c(\theta, a, \gamma) v(\gamma) \leq \boldsymbol{\rho}_t, \\ &\sum_{a \in A^+} \phi(\theta, a) \leq 1, \quad \forall \theta \in \Theta, \\ &\phi(\theta, a) \geq 0, \quad \forall (\theta, a) \in \Theta \times A^+. \end{aligned}$$

$$\begin{aligned} \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) &:= \max_{\phi: \Theta \times A^+ \rightarrow \mathbb{R}_+} \sum_{\theta \in \Theta, a \in A^+} \hat{u}_t(\theta) \phi(\theta, a) \sum_{\gamma} r(\theta, a, \gamma) \hat{v}_t(\gamma), \\ \text{s.t.} \quad &\sum_{\theta \in \Theta, a \in A^+} \hat{u}_t(\theta) \phi(\theta, a) \sum_{\gamma} c(\theta, a, \gamma) \hat{v}_t(\gamma) \leq \boldsymbol{\rho}_t, \\ &\sum_{a \in A^+} \phi(\theta, a) \leq 1, \quad \forall \theta \in \Theta, \\ &\phi(\theta, a) \geq 0, \quad \forall (\theta, a) \in \Theta \times A^+. \end{aligned}$$

747 Now, conceptually, if there is a $0 < \eta_t \leq 1$ such that for any $(\theta, a) \in \Theta \times A^+$,

$$u(\theta) \sum_{\gamma} c(\theta, a, \gamma) v(\gamma) \geq \eta_t \hat{u}_t(\theta) \sum_{\gamma} c(\theta, a, \gamma) \hat{v}_t(\gamma),$$

748 then for an optimal solution ϕ_t^* of $J(\boldsymbol{\rho}_t)$, we see that $\eta_t \phi_t^*$ is a feasible solution of the programming
 749 $\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)$. Thus,

$$\begin{aligned} \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) &\geq \eta_t \sum_{\theta \in \Theta, a \in A^+} \hat{u}_t(\theta) \phi_t^*(\theta, a) \sum_{\gamma} r(\theta, a, \gamma) \hat{v}_t(\gamma) \\ &\stackrel{(a)}{\geq} \eta_t \sum_{\theta \in \Theta, a \in A^+} u(\theta) \phi_t^*(\theta, a) \sum_{\gamma} r(\theta, a, \gamma) v(\gamma) \\ &\quad - \eta_t (\| (u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta} \|_1 + \| (v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma} \|_1) \\ &= \eta_t J(\boldsymbol{\rho}_t) - (\| (u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta} \|_1 + \| (v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma} \|_1). \end{aligned}$$

750 Here, since $r(\theta, a, \gamma) \leq 1$ and $\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \leq 1$ for any θ , (a) is expanded as

$$\begin{aligned} &\sum_{\theta \in \Theta, a \in A^+} \hat{u}_t(\theta) \phi_t^*(\theta, a) \sum_{\gamma} \hat{v}_t(\gamma) r(\theta, a, \gamma) - \sum_{\theta \in \Theta, a \in A^+} u(\theta) \phi_t^*(\theta, a) \sum_{\gamma} v(\gamma) r(\theta, a, \gamma) \\ &= \sum_{\theta \in \Theta, a \in A^+} (\hat{u}_t(\theta) - u(\theta)) \phi_t^*(\theta, a) \sum_{\gamma} \hat{v}_t(\gamma) r(\theta, a, \gamma) \\ &\quad + \sum_{\theta \in \Theta, a \in A^+} u(\theta) \phi_t^*(\theta, a) \sum_{\gamma} (\hat{v}_t(\gamma) - v(\gamma)) r(\theta, a, \gamma) \\ &\leq \| (u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta} \|_1 + \| (v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma} \|_1. \end{aligned}$$

751 Consequently,

$$\begin{aligned} & J(\boldsymbol{\rho}_1) - \widehat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) \\ & \leq (1 - \eta_t)J(\boldsymbol{\rho}_1) + \eta_t(J(\boldsymbol{\rho}_1) - J(\boldsymbol{\rho}_t)) + \|(u(\theta) - \widehat{u}_t(\theta))_{\theta \in \Theta}\|_1 + \|(v(\gamma) - \widehat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1. \end{aligned} \quad (25)$$

752 On top of this, a key observation is that

$$J(\boldsymbol{\rho}_1) - J(\boldsymbol{\rho}_t) \leq \frac{\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)}{\rho^{\min}} \cdot J(\boldsymbol{\rho}_1). \quad (26)$$

753 In fact, when $\boldsymbol{\rho}_1 \leq \boldsymbol{\rho}_t$, (26) is natural as $J(\boldsymbol{\rho}_1) \leq J(\boldsymbol{\rho}_t)$. Otherwise, let ϕ_1^* be the optimal solution to
754 the programming $J(\boldsymbol{\rho}_1)$. Let i^* be the index that minimizes $\rho_t^{i^*}/\rho_1^{i^*}$. We have $\rho_1^{i^*} > \rho_t^{i^*}$. Evidently,
755 we know that $\phi_1^* \cdot \rho_t^{i^*}/\rho_1^{i^*}$ is a feasible solution to the programming of $J(\boldsymbol{\rho}_t)$. By the optimality of
756 $J(\boldsymbol{\rho}_t)$, we have

$$J(\boldsymbol{\rho}_t) \geq \frac{\rho_t^{i^*}}{\rho_1^{i^*}} \cdot J(\boldsymbol{\rho}_1),$$

757 which leads to

$$J(\boldsymbol{\rho}_1) - J(\boldsymbol{\rho}_t) \leq \left(1 - \frac{\rho_t^{i^*}}{\rho_1^{i^*}}\right) \cdot J(\boldsymbol{\rho}_1) = \frac{\rho_1^{i^*} - \rho_t^{i^*}}{\rho_1^{i^*}} \cdot J(\boldsymbol{\rho}_1) \leq \frac{\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t)}{\rho^{\min}} \cdot J(\boldsymbol{\rho}_1).$$

758 Synthesizing the above two parts, (26) is proved.

759 As for $\mathbb{E}[\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)]$, we note that for any non-negative random variable X with upper bound
760 \bar{X} and any positive ξ , we have

$$\mathbb{E}[X] \leq \xi \Pr[X \leq \xi] + \bar{X} (1 - \Pr[X \leq \xi]) \leq \xi + \bar{X} (1 - \Pr[X \leq \xi]). \quad (27)$$

Notice that $\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)$ is certainly upper bounded by 1. Therefore, as a corollary of Lemma E.1, we have

$$\mathbb{E}[\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)] \leq \begin{cases} O\left(\frac{\sqrt{(t-2)\log T}}{T} + \sqrt{\frac{\log T}{T-t}} + \frac{1}{T}\right), & 2 \leq t \leq (T+1)/2; \\ O\left(\sqrt{\frac{\log T}{T-t}} + \frac{1}{T}\right), & t > (T+1)/2. \end{cases}$$

761 We almost finish the bound now except for determining η_t in (25), which we hope is as close to 1 as
762 possible. Nevertheless, we leave the technical parts to Appendix E.5 which derives the following
763 lemma on the total bound:

Lemma E.3.

$$\mathbb{E}\left[\sum_{t=1}^{T_0} \left(J(\boldsymbol{\rho}_1) - \widehat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)\right)\right] = O(\sqrt{T \log T}).$$

764 Now, we sum the result in (21) from $t = 2$ to T_0 , and plus the constant term for $t = 1$ to obtain that

$$\mathbb{E}\left[\sum_{t=1}^{T_0} \left(\widehat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_{\theta} \left[\sum_{a \in A^+} \widehat{\phi}_t^*(\theta, a) R(\theta, a)\right]\right)\right] = O(\sqrt{T \log T}).$$

765 Synthesizing Lemma E.3, (24), (23), and (21), we derive Theorem 5.1.

766 E.2 Proof of Theorem 5.2

767 The proof of this theorem follows the line of Theorem 5.1, and the only difference is to adopt
768 Lemma 4.1 when considering the concentration of estimates. On this side, we can disregard the cases
769 when $t \leq \Theta(\log T)$, as the accumulated regret in this phase is bounded by $O(\log T)$. On the other
770 hand, the time range that $t \geq \Theta(T)$ is asymptotically identical to the full information setting since

the accessing frequency is a constant. We only need to consider the case that $\Theta(\log T) \leq t \leq \Theta(T)$, when we have

$$\begin{aligned} \Pr \left[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1 \leq -\Theta \left(\sqrt{\frac{\log T}{t-1}} \right) \right] &\leq O \left(\frac{1}{T^2} \right), \\ \Pr \left[\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \leq -\Theta \left(\frac{\log T}{\sqrt{t-1}} \right) \right] &\leq O \left(\frac{1}{T^2} \right). \end{aligned} \quad (28)$$

Taking into the proof of Lemma E.1 and then into the main body, we should find a sufficient large t such that

$$\begin{aligned} &\Theta \left(\frac{\log T}{T-1} + \sum_{\tau=\Theta(\log T)}^{\Theta(T)} \frac{\log T}{\sqrt{(T-\tau)^2(\tau-1)}} + \sum_{\tau=\Theta(T)}^{t-1} \sqrt{\frac{\log T}{(T-\tau)^2(\tau-1)}} + \sqrt{\frac{\log T}{T-t}} \right) \\ &\leq \rho^{\min} - \frac{1}{T-t+1}, \end{aligned}$$

and $t = T - O(\log T)$ still suffices. Therefore, $\mathbb{E}[T - T_0] = O(\log T)$ also holds under partial information feedback.

Nevertheless, for the counterpart of Lemma E.2, by (28), when we sum from $t = 1$ to $T_0 \leq T$, we derive that

$$\mathbb{E} \left[\sum_{t=1}^{T_0} \left(\hat{J}(\rho_t, \mathcal{H}_t) - \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \right) \right] \leq O(\sqrt{T} \log T).$$

At last, for $J(\rho_1) - \hat{J}(\rho_t, \mathcal{H}_t)$, we face the same degradation on the estimation accuracy, which leads to

$$\mathbb{E} \left[\sum_{t=1}^{T_0} \left(J(\rho_1) - \hat{J}(\rho_t, \mathcal{H}_t) \right) \right] = O(\sqrt{T} \log T).$$

Therefore, Theorem 5.2 is achieved.

E.3 Proof of Lemma E.1

Now that we are going to bound $\max(\rho_1 - \rho_t, 0)$. Recall the definitions below which we give in Appendix C.5 when we prove Lemma C.4:

$$\mathbf{M}_t^C := \frac{\rho_t - \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \mathbf{C}(\theta, a) \right]}{T-t}, \quad \mathbf{N}_t^C := \frac{\mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \mathbf{C}(\theta, a) \right] - \mathbf{c}_t}{T-t}.$$

By (16), we have

$$\rho_{t+1} - \rho_t = \frac{\rho_t - \mathbf{c}_t}{T-t} = \mathbf{M}_t^C + \mathbf{N}_t^C.$$

Consequently,

$$\max(\rho_1 - \rho_t) = \max \left(- \left(\sum_{\tau=1}^{t-1} \mathbf{M}_\tau^C + \sum_{\tau=1}^{t-1} \mathbf{N}_\tau^C \right) \right) \leq - \min \sum_{\tau=1}^{t-1} \mathbf{M}_\tau^C - \min \sum_{\tau=1}^{t-1} \mathbf{N}_\tau^C.$$

For the second term, we notice that each entry of $\{\sum_{\tau < t} \mathbf{N}_\tau^C\}_t$ is a martingale with the absolute value of the τ -th increment bounded by $1/(T-\tau)$. Since

$$\sum_{\tau=1}^{t-1} \frac{1}{(T-\tau)^2} \leq \frac{1}{T-t},$$

by applying the Azuma–Hoeffding inequality and a union bound, we achieve that

$$\Pr \left[- \min \sum_{\tau=1}^{t-1} \mathbf{N}_\tau^C \geq \sqrt{\frac{2 \log T}{T-t}} \right] \leq \frac{n}{T}. \quad (29)$$

On the other hand, for the first term, when $\tau = 1$, it is apparent that $-\min \mathbf{M}_1^C \leq 1/(T-1)$. When $\tau \geq 2$, we have for any $i \in [n]$,

$$\begin{aligned}
& (T - \tau) (\mathbf{M}_\tau^C)^i \\
&= \boldsymbol{\rho}_\tau^i - \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \mathbf{C}^i(\theta, a) \right] \\
&\stackrel{(a)}{\geq} \sum_{\theta \in \Theta, a \in A^+} \hat{u}_\tau(\theta) \hat{\phi}_\tau^*(\theta, a) \sum_{\gamma} \hat{v}_\tau(\gamma) \mathbf{c}^i(\theta, a, \gamma) - \sum_{\theta \in \Theta, a \in A^+} u(\theta) \hat{\phi}_\tau^*(\theta, a) \sum_{\gamma} v(\gamma) \mathbf{c}^i(\theta, a, \gamma) \\
&= \sum_{\theta \in \Theta, a \in A^+} (\hat{u}_\tau(\theta) - u(\theta)) \hat{\phi}_\tau^*(\theta, a) \sum_{\gamma} \hat{v}_\tau(\gamma) \mathbf{c}^i(\theta, a, \gamma) \\
&\quad + \sum_{\theta \in \Theta, a \in A^+} u(\theta) \hat{\phi}_\tau^*(\theta, a) \sum_{\gamma} (\hat{v}_\tau(\gamma) - v(\gamma)) \mathbf{c}^i(\theta, a, \gamma) \\
&\geq -\|(u(\theta) - \hat{u}_\tau(\theta))_{\theta \in \Theta}\|_1 - \|(v(\gamma) - \hat{v}_\tau(\gamma))_{\gamma \in \Gamma}\|_1.
\end{aligned}$$

In the above, (a) is because $\hat{\phi}_\tau^*$ is feasible for $\hat{J}(\boldsymbol{\rho}_\tau, \mathcal{H}_\tau)$. By Hoeffding's inequality and a union bound, we have

$$\begin{aligned}
& \Pr \left[\|(u(\theta) - \hat{u}_\tau(\theta))_{\theta \in \Theta}\|_1 \leq -|\Theta| \sqrt{\frac{\log T}{\tau - 1}} \right] \leq \frac{|\Theta|}{T^2}, \\
& \Pr \left[\|(v(\gamma) - \hat{v}_\tau(\gamma))_{\gamma \in \Gamma}\|_1 \leq -|\Gamma| \sqrt{\frac{\log T}{\tau - 1}} \right] \leq \frac{|\Gamma|}{T^2}.
\end{aligned}$$

Thus, suppose the above events hold for all $\tau \leq T$ with failure probability only $O(1/T)$,

$$\Pr \left[-\min_{\tau=1}^{t-1} \mathbf{M}_\tau^C \geq \Theta \left(\frac{1}{T-1} + \sum_{\tau=2}^{t-1} \sqrt{\frac{\log T}{(T-\tau)^2(\tau-1)}} \right) \right] \leq O\left(\frac{1}{T}\right). \quad (30)$$

Combining (29) and (30), we derive the lemma.

E.4 Proof of Lemma E.2

We notice that

$$\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) = \sum_{\theta \in \Theta, a \in A^+} \hat{u}_t(\theta) \hat{\phi}_t^*(\theta, a) \sum_{\gamma} r(\theta, a, \gamma) \hat{v}_t(\gamma),$$

and

$$\begin{aligned}
& \sum_{\theta \in \Theta, a \in A^+} \hat{u}_t(\theta) \hat{\phi}_t^*(\theta, a) \sum_{\gamma} \hat{v}_t(\gamma) r(\theta, a, \gamma) - \sum_{\theta \in \Theta, a \in A^+} u(\theta) \hat{\phi}_t^*(\theta, a) \sum_{\gamma} v(\gamma) r(\theta, a, \gamma) \\
&= \sum_{\theta \in \Theta, a \in A^+} (\hat{u}_t(\theta) - u(\theta)) \hat{\phi}_t^*(\theta, a) \sum_{\gamma} \hat{v}_t(\gamma) r(\theta, a, \gamma) \\
&\quad + \sum_{\theta \in \Theta, a \in A^+} u(\theta) \hat{\phi}_t^*(\theta, a) \sum_{\gamma} (\hat{v}_t(\gamma) - v(\gamma)) r(\theta, a, \gamma) \\
&\leq \|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1 + \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1.
\end{aligned}$$

Thus,

$$\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_{\theta \sim \mathcal{U}} \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \leq \|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1 + \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1.$$

800 By Hoeffding's inequality and a union bound, we have

$$\begin{aligned} \Pr \left[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1 \geq |\Theta| \sqrt{\frac{\log T}{2(t-1)}} \right] &\leq \frac{|\Theta|}{T}, \\ \Pr \left[\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \geq |\Gamma| \sqrt{\frac{\log T}{2(t-1)}} \right] &\leq \frac{|\Gamma|}{T}. \end{aligned}$$

801 Further, the difference we hope to analyze is certainly upper bounded by 1. As a result, with (27), we
802 finish the proof.

803 E.5 Proof of Lemma E.3

804 We come to consider $J(\boldsymbol{\rho}_1) - \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)$. As per the thread in the main body, we let

$$\delta_t := \frac{\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_\infty + \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1}{\min_{\theta \in \Theta, a \in A^+} \{\min\{u(\theta) \mathbf{C}(\theta, a) > 0\}\}}.$$

805 We now claim that for any $(\theta, a, i) \in \Theta \times A^+ \times [n]$,

$$\hat{u}_t(\theta) \sum_{\gamma} \mathbf{c}^i(\theta, a, \gamma) \hat{v}_t(\gamma) \leq (1 + \delta_t) u(\theta) \sum_{\gamma} \mathbf{c}^i(\theta, a, \gamma) v(\gamma).$$

806 The above is obvious if $\mathbf{C}^i(\theta, a) = \mathbf{0}$, or $\mathbf{c}^i(\theta, a, \gamma) = 0$ holds for any γ . When $\mathbf{C}(\theta, a) \neq \mathbf{0}$, then
807 for any $i \in [n]$,

$$\begin{aligned} &\hat{u}_t(\theta) \sum_{\gamma} \mathbf{c}^i(\theta, a, \gamma) \hat{v}_t(\gamma) - u(\theta) \sum_{\gamma} \mathbf{c}^i(\theta, a, \gamma) v(\gamma) \\ &= (\hat{u}_t(\theta) - u(\theta)) \sum_{\gamma} \mathbf{c}^i(\theta, a, \gamma) \hat{v}_t(\gamma) + u(\theta) \sum_{\gamma} \mathbf{c}^i(\theta, a, \gamma) (\hat{v}_t(\gamma) - v(\gamma)) \\ &\leq \|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_\infty + \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \\ &\leq \delta_t u(\theta) \sum_{\gamma} \mathbf{c}^i(\theta, a, \gamma) v(\gamma). \end{aligned}$$

808 This finish the explanation of the claim. Upon that, if we let $\eta_t := 1 - \delta_t \leq 1/(1 + \delta_t)$, we derive that

$$\begin{aligned} u(\theta) \sum_{\gamma} \mathbf{c}(\theta, a, \gamma) v(\gamma) &\leq \frac{1}{1 + \delta_t} \hat{u}_t(\theta) \sum_{\gamma} \mathbf{c}(\theta, a, \gamma) \hat{v}_t(\gamma) \\ &\leq \eta_t \hat{u}_t(\theta) \sum_{\gamma} \mathbf{c}(\theta, a, \gamma) \hat{v}_t(\gamma). \end{aligned}$$

809 With respect to (25) and (26), we obtain that

$$\begin{aligned} &J(\boldsymbol{\rho}_1) - \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) \\ &\leq J(\boldsymbol{\rho}_1) \cdot \left(1 - \eta_t + \frac{\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)}{\rho^{\min}} \right) + \|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1 + \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \\ &= J(\boldsymbol{\rho}_1) \cdot \left(\delta_t + \frac{\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)}{\rho^{\min}} \right) + \|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1 + \|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1. \quad (31) \end{aligned}$$

As we have already shown in the main body that

$$\mathbb{E}[\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)] \leq \begin{cases} O\left(\frac{\sqrt{(t-2)\log T}}{T} + \sqrt{\frac{\log T}{T-t}} + \frac{1}{T}\right), & 2 \leq t \leq (T+1)/2; \\ O\left(\sqrt{\frac{\log T}{T-t}} + \frac{1}{T}\right), & t > (T+1)/2, \end{cases}$$

810 it suffices for us to bound

$$\mathbb{E}[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_\infty], \mathbb{E}[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1], \mathbb{E}[\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1].$$

811 On this side, as we have shown that

$$\begin{aligned} \Pr \left[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1 \geq |\Theta| \sqrt{\frac{\log T}{2(t-1)}} \right] &\leq \frac{|\Theta|}{T}, \\ \Pr \left[\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1 \geq |\Gamma| \sqrt{\frac{\log T}{2(t-1)}} \right] &\leq \frac{|\Gamma|}{T}, \end{aligned}$$

812 it is natural that

$$\begin{aligned} &\{\mathbb{E}[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_\infty], \mathbb{E}[\|(u(\theta) - \hat{u}_t(\theta))_{\theta \in \Theta}\|_1], \mathbb{E}[\|(v(\gamma) - \hat{v}_t(\gamma))_{\gamma \in \Gamma}\|_1]\} \\ &\leq O \left(\sqrt{\frac{\log T}{t-1}} + \frac{1}{T} \right). \end{aligned}$$

813 Thus, putting all the above into (31) and summing from $t = 1$ to $T_0 \leq T$, we have

$$\mathbb{E} \left[\sum_{t=1}^{T_0} \left(J(\rho_1) - \hat{J}(\rho_t, \mathcal{H}_t) \right) \right] = O(\sqrt{T \log T}).$$

814 This concludes the proof.

815 **F Missing Details in Appendix A**

816 **F.1 The Density Estimator**

817 We now present details on the kernel density estimator which we apply in Appendix A for approximat-
818 ing continuous distributions, which comes from Wasserman [2019]. We consider a one-dimensional
819 kernel function K such that

- 820 • $\int K(x) dx = 1$;
- 821 • $\int x^s K(x) dx = 0, \quad \forall 1 \leq s \leq \beta$;
- 822 • $\int |x|^\beta |K(x)| dx < \infty$.

823 Now, given k independent samples X_1, \dots, X_k from P and a positive number h called the bandwidth,
824 the kernel density estimator is defined as

$$\hat{p}_k(x) = \frac{1}{k} \sum_{i=1}^k \frac{1}{h^d} K \left(\frac{\|x - X_i\|_2}{h} \right).$$

825 Furthermore, to satisfy Proposition A.1, we should choose $h \asymp k^{1/(2\beta+d)} \log k$ when $p \in \Sigma(\beta, L)$ is
826 the density of \mathcal{P} on \mathbb{R}^d .

827 **F.2 Proof of Theorem A.1**

828 By (21), we know that

$$V^{\text{FL}} - \text{Rew} = J(\rho_1) \cdot \mathbb{E}[T - T_0] + \mathbb{E} \left[\sum_{t=1}^{T_0} \left(J(\rho_1) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \right) \right],$$

829 and we bound these terms in order. For the expected stopping time $\mathbb{E}[T_0]$, by the analysis in Section 5,
830 our goal turns into bounding $\max(\rho_1 - \rho_t, 0)$, which further by (16) and (29), reduces to bound

831 M_τ^C . With continuous randomness, we have for any $i \in [n]$,

$$\begin{aligned}
& (T - \tau) (M_\tau^C)^i \\
&= \rho_\tau^i - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \mathbf{C}^i(\theta, a) \right] \\
&\stackrel{(a)}{\geq} \int_\theta \sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) \hat{v}_\tau(\gamma) \hat{u}_\tau(\theta) d\gamma d\theta \\
&\quad - \int_\theta \sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) v(\gamma) u(\theta) d\gamma d\theta \\
&= \int_\theta \sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) \hat{v}_\tau(\gamma) (\hat{u}_\tau(\theta) - u(\theta)) d\gamma d\theta \\
&\quad + \int_\theta \sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) (\hat{v}_\tau(\gamma) - v(\gamma)) u(\theta) d\gamma d\theta \\
&\stackrel{(b)}{\geq} - \sup_\theta |u(\theta) - \hat{u}_\tau(\theta)| - \sup_\gamma |v(\gamma) - \hat{v}_\tau(\gamma)|.
\end{aligned}$$

832 In the above, (a) is by the constraint feasibility of $\hat{\phi}_\tau^*$, and (b) is because $\sum_{a \in A^+} \hat{\phi}_\tau^*(\theta, a) \leq 1$ holds
833 for any $\theta \in \Theta$. Further, by Proposition A.1, we have for $\tau = \Omega(1)$,

$$\begin{aligned}
& \Pr \left[\sup_\theta |u(\theta) - \hat{u}_\tau(\theta)| \leq -\Theta \left(\sqrt{\log T} (\tau - 1)^{\alpha_u - 1} \right) \right] \leq \frac{1}{T^2}, \\
& \Pr \left[\sup_\gamma |v(\gamma) - \hat{v}_\tau(\gamma)| \leq -\Theta \left(\sqrt{\log T} (\tau - 1)^{\alpha_v - 1} \right) \right] \leq \frac{1}{T^2}.
\end{aligned}$$

834 Thus, when $t = \Omega(1)$, we derive that with failure probability $O(1/T)$, it holds that

$$\max(\rho_1 - \rho_t, 0) \leq \Theta \left(\frac{1}{T-1} + \sqrt{\log T} \sum_{\tau=\Theta(1)}^{t-1} \left(\frac{(\tau-1)^{\alpha_u-1}}{T-\tau} + \frac{(\tau-1)^{\alpha_v-1}}{T-\tau} \right) + \sqrt{\frac{\log T}{T-t}} \right).$$

835 Further, for $p \in \{u, v\}$, when $t \leq (T+1)/2$,

$$\sum_{\tau=\Theta(1)}^{t-1} \frac{(\tau-1)^{\alpha_p-1}}{T-\tau} \leq \frac{2}{T-1} \sum_{\tau=2}^{t-1} (\tau-1)^{\alpha_p-1} \leq \frac{2(t-2)^{\alpha_p}}{\alpha_p(T-1)};$$

836 and when $t > (T+1)/2$, we have

$$\sum_{\tau=\Theta(1)}^{t-1} \frac{(\tau-1)^{\alpha_p-1}}{T-\tau} \leq \frac{1}{\alpha_p} \left(\frac{2}{T-1} \right)^{1-\alpha_p} + \sum_{\tau=(T+1)/2}^{t-1} (T-\tau)^{\alpha_p-2} \leq \frac{(T-t)^{\alpha_p-1}}{1-\alpha_p}.$$

837 Thus, when $t = T - \Theta(\log^{(2(1-\max\{1/2, \alpha_u, \alpha_v\}))^{-1}} T)$, we have

$$\begin{aligned}
& \Theta \left(\frac{1}{T-1} + \sqrt{\log T} \sum_{\tau=\Theta(1)}^{t-1} \left(\frac{(\tau-1)^{\alpha_u-1}}{T-\tau} + \frac{(\tau-1)^{\alpha_v-1}}{T-\tau} \right) + \sqrt{\frac{\log T}{T-t}} \right) \\
& \leq \rho^{\min} - \frac{1}{T-t+1},
\end{aligned}$$

838 which leads to

$$\mathbb{E}[T - T_0] = O \left(\log^{(2(1-\max\{1/2, \alpha_u, \alpha_v\}))^{-1}} T \right).$$

839 This concludes the analysis of the stopping time.

840 For the second part, By (24), we have

$$\begin{aligned} & J(\boldsymbol{\rho}_1) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \\ &= \left(J(\boldsymbol{\rho}_1) - \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) \right) + \left(\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \right). \end{aligned}$$

841 On the second difference term, similar to the proof of Lemma E.2, we have

$$\begin{aligned} & \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \\ &= \int_\theta \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \int_\gamma r(\theta, a, \gamma) \hat{v}_t(\gamma) \hat{u}_t(\theta) d\gamma d\theta - \int_\theta \sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) \int_\gamma r(\theta, a, \gamma) v(\gamma) u(\theta) d\gamma d\theta \\ &\leq \sup_\theta |u(\theta) - \hat{u}_t(\theta)| + \sup_\gamma |(v(\gamma) - \hat{v}_t(\gamma))|. \end{aligned}$$

842 Thus, when $t = \Omega(1)$, by taking $\epsilon = 1/T$ in Proposition A.1 and (27), we arrive that

$$\mathbb{E} \left[\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_\theta \left[\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a) \right] \right] = O \left(\sqrt{\log T} ((t-1)^{\alpha_u-1} + (t-1)^{\alpha_v-1}) + \frac{1}{T} \right).$$

843 We now focus on $J(\boldsymbol{\rho}_1) - \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)$. We let

$$\delta_t := \frac{\sup_\theta |u(\theta) - \hat{u}_t(\theta)| + \sup_\gamma |(v(\gamma) - \hat{v}_t(\gamma))|}{\min_{\theta \in \Theta, a \in A^+} \{\min\{u(\theta)C(\theta, a) > 0\}\}}.$$

844 We prove that

$$\hat{u}_t(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) \hat{v}_t(\gamma) d\gamma \leq (1 + \delta_t) u(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) v(\gamma) d\gamma$$

845 holds for any (θ, a, i) tuple, which is obvious if $\mathbf{c}^i(\theta, a, \gamma)$ is almost surely zero with respect to γ .

846 Otherwise, we observe that

$$\begin{aligned} & \hat{u}_t(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) \hat{v}_t(\gamma) d\gamma - u(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) v(\gamma) d\gamma \\ &= (\hat{u}_t(\theta) - u(\theta)) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) \hat{v}_t(\gamma) d\gamma + u(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) (\hat{v}_t(\gamma) - v(\gamma)) d\gamma \\ &\leq \sup_\theta |u(\theta) - \hat{u}_t(\theta)| + \sup_\gamma |(v(\gamma) - \hat{v}_t(\gamma))| \\ &\leq \delta_t u(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) v(\gamma) d\gamma. \end{aligned}$$

847 and thus, with $\eta_t := 1 - \delta_t \leq 1/(1 + \delta_t)$, we derive that

$$\begin{aligned} u(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) v(\gamma) d\gamma &\leq \frac{1}{1 + \delta_t} \hat{u}_t(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) \hat{v}_t(\gamma) d\gamma \\ &\leq \eta_t \hat{u}_t(\theta) \int_\gamma \mathbf{c}^i(\theta, a, \gamma) \hat{v}_t(\gamma) d\gamma. \end{aligned}$$

848 This proves the above inequality. Thus, for an optimal solution ϕ_t^* of $J(\boldsymbol{\rho}_t)$, we see that $\eta_t \phi_t^*$ is a
 849 feasible solution of the programming $\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)$. Thus, we notice that

$$\begin{aligned}\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) &\geq \eta_t \int_{\theta} \sum_{a \in A^+} \phi_t^*(\theta, a) \int_{\gamma} r(\theta, a, \gamma) \hat{v}_t(\gamma) \hat{u}_t(\theta) d\gamma d\theta \\ &\geq \eta_t \int_{\theta} \sum_{a \in A^+} \phi_t^*(\theta, a) \int_{\gamma} r(\theta, a, \gamma) v(\gamma) u(\theta) d\gamma d\theta \\ &\quad - \eta_t (\sup_{\theta} |u(\theta) - \hat{u}_t(\theta)| + \sup_{\gamma} |(v(\gamma) - \hat{v}_t(\gamma))|) \\ &= \eta_t J(\boldsymbol{\rho}_t) - (\sup_{\theta} |u(\theta) - \hat{u}_t(\theta)| + \sup_{\gamma} |(v(\gamma) - \hat{v}_t(\gamma))|).\end{aligned}$$

850 With respect to (26), we obtain that

$$\begin{aligned}&J(\boldsymbol{\rho}_1) - \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) \\ &\leq J(\boldsymbol{\rho}_1) \cdot \left(1 - \eta_t + \frac{\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)}{\rho^{\min}}\right) + \sup_{\theta} |u(\theta) - \hat{u}_t(\theta)| + \sup_{\gamma} |(v(\gamma) - \hat{v}_t(\gamma))| \\ &= J(\boldsymbol{\rho}_1) \cdot \left(\delta_t + \frac{\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)}{\rho^{\min}}\right) + \sup_{\theta} |u(\theta) - \hat{u}_t(\theta)| + \sup_{\gamma} |(v(\gamma) - \hat{v}_t(\gamma))|.\end{aligned}$$

851 Now, when $t = \Theta(1)$, we have

$$\begin{aligned}\mathbb{E} \left[\sup_{\theta} |u(\theta) - \hat{u}_t(\theta)| \right] &= O \left(\sqrt{\log T} (t-1)^{\alpha_u-1} + \frac{1}{T} \right), \\ \mathbb{E} \left[\sup_{\gamma} |v(\gamma) - \hat{v}_t(\gamma)| \right] &= O \left(\sqrt{\log T} (t-1)^{\alpha_v-1} + \frac{1}{T} \right).\end{aligned}$$

852 By the previous reasoning on $\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)$, we obtain that when $t = \Omega(1)$,

$$\begin{aligned}&\mathbb{E} [\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)] \\ &\leq \Theta \left(\frac{1}{T-1} + \sqrt{\log T} \sum_{\tau=\Theta(1)}^{t-1} \left(\frac{(\tau-1)^{\alpha_u-1}}{T-\tau} + \frac{(\tau-1)^{\alpha_v-1}}{T-\tau} \right) + \sqrt{\frac{\log T}{T-t}} \right).\end{aligned}$$

853 Therefore, summing from $t = 1$ to $T_0 \leq T$, we achieve that

$$\mathbb{E} \left[\sum_{t=1}^{T_0} \left(J(\boldsymbol{\rho}_1) - \hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) \right) \right] = O \left((T^{1/2} + T^{\alpha_u} + T^{\alpha_v}) \sqrt{\log T} \right).$$

854 Combining with previous bounds on $\mathbb{E}[T - T_0]$ and the estimation errors, we derive the theorem.

855 F.3 Proof of Theorem A.2

856 Similar to the proof of Theorem 5.2, we concentrate on re-bounding the three terms under partial
 857 information feedback, respectively $\mathbb{E}[T - T_0]$, $\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_{\theta} [\sum_{a \in A^+} \hat{\phi}_t^*(\theta, a) R(\theta, a)]$, and $J(\boldsymbol{\rho}_1) -$
 858 $\hat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)$. As for $\mathbb{E}[T - T_0]$, with Lemma 4.1, we argue here that the main term in bounding
 859 $\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)$ when $t = \Theta(T)$ becomes

$$\Theta \left(\sqrt{\log T} \left(\sum_{\tau=\Theta(1)}^{t-1} \frac{(\tau-1)^{\alpha_u-1}}{T-\tau} + \sum_{\tau=\Theta(1)}^{\Theta(T)} \frac{((\tau-1)/\log T)^{\alpha_v-1}}{T-\tau} + \sum_{\tau=\Theta(T)}^{t-1} \frac{(\tau-1)^{\alpha_v-1}}{T-\tau} \right) \right).$$

860 Consequently, when t is close to T , we have with failure probability $O(1/T)$,

$$\begin{aligned}&\max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0) \\ &\leq \Theta \left(\frac{1}{T-1} + \sqrt{\log T} \left((T-t)^{\alpha_u-1} + (T-t)^{-1/2} \right) + (T-t)^{\alpha_v-1} \log^{3/2-\alpha_v} T \right).\end{aligned}$$

861 This leads to

$$\mathbb{E}[T - T_0] = O\left(\log^{\max(1, 1/(2-2\alpha_u), (3-2\alpha_v)/(2-2\alpha_v))} T\right).$$

862 For the estimation error term $\widehat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_\theta[\sum_{a \in A^+} \widehat{\phi}_t^*(\theta, a)R(\theta, a)]$, when $\Omega(1) \leq t \leq \Theta(T)$,
 863 the bound now becomes

$$\begin{aligned} & \mathbb{E}\left[\widehat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t) - \mathbb{E}_\theta\left[\sum_{a \in A^+} \widehat{\phi}_t^*(\theta, a)R(\theta, a)\right]\right] \\ &= O\left(\sqrt{\log T}(t-1)^{\alpha_u-1} + \log^{3/2-\alpha_v} T \cdot (t-1)^{\alpha_v-1} + \frac{1}{T}\right). \end{aligned}$$

864 At last, for $J(\boldsymbol{\rho}_1) - \widehat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)$, we derive that

$$\begin{aligned} & \mathbb{E}\left[\sum_{t=1}^{T_0} \max(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_t, 0)\right] = O\left(\sqrt{\log T} \left(T^{\alpha_u} + T^{1/2}\right) + \log^{3/2-\alpha_v} T \cdot T^{\alpha_v}\right), \\ & \mathbb{E}\left[\sum_{t=1}^{T_0} \sup_{\theta} |u(\theta) - \widehat{u}_t(\theta)|\right] = O\left(\sqrt{\log T} \cdot T^{\alpha_u}\right), \\ & \mathbb{E}\left[\sum_{t=1}^{T_0} \sup_{\gamma} |v(\gamma) - \widehat{v}_t(\gamma)|\right] = O\left(\log^{3/2-\alpha_v} T \cdot T^{\alpha_v}\right). \end{aligned}$$

865 Putting together, we obtain that

$$\mathbb{E}\left[\sum_{t=1}^{T_0} \left(J(\boldsymbol{\rho}_1) - \widehat{J}(\boldsymbol{\rho}_t, \mathcal{H}_t)\right)\right] = O\left(\sqrt{\log T} \left(T^{\alpha_u} + T^{1/2}\right) + \log^{3/2-\alpha_v} T \cdot T^{\alpha_v}\right).$$

866 Synthesizing all the above, we finish the proof of the theorem.