
Reinforcement Learning with History-Dependent Dynamic Contexts

Guy Tennenholtz^{*1} Nadav Merlis^{*2} Lior Shani¹ Martin Mladenov¹ Craig Boutilier¹

Abstract

We introduce *Dynamic Contextual Markov Decision Processes (DCMDPs)*, a novel reinforcement learning framework for history-dependent environments that generalizes the contextual MDP framework to handle non-Markov environments, where contexts change over time. We consider special cases of the model, with a focus on *logistic DCMDPs*, which break the exponential dependence on history length by leveraging aggregation functions to determine context transitions. This special structure allows us to derive an upper-confidence-bound style algorithm for which we establish regret bounds. Motivated by our theoretical results, we introduce a practical model-based algorithm for logistic DCMDPs that plans in a latent space and uses optimism over history-dependent features. We demonstrate the efficacy of our approach on a recommendation task (using MovieLens data) where user behavior dynamics evolve in response to recommendations.

1. Introduction

Reinforcement learning (RL) is a paradigm in which an agent learns to act in an environment to maximize long-term reward. RL has been applied to numerous domains, including recommender systems, robot control, video games, and autonomous vehicles (Afsar et al., 2022; Tessler et al., 2019; Mnih et al., 2015; Fayjie et al., 2018). While typical RL approaches rely on a Markov property of both the reward process and environment dynamics, many scenarios are inherently history-dependent (Bacchus et al., 1996; Ronca and Giacomo, 2021), particularly, when humans are involved. As one example, the behavior of users in recommender systems often exhibits non-Markovian characteristics reflective of a user’s latent state, including: user preference elicitation sessions, where users respond to a sequence of

feedback-gathering interventions (e.g., ratings, comparisons, annotations) (Chen and Pu, 2012; Zhao et al., 2013); user *ad blindness* (i.e., the tendency to gradually ignore ads) (Hohnhold et al., 2015); and the long-term evolution of user satisfaction (Wilhelm et al., 2018; Mladenov et al., 2019). Many aspects of a user’s latent state determine their disposition towards specific actions. For example, a user’s level of frustration, trust, receptivity, and overall satisfaction, may affect their tendency toward accepting recommendations, providing feedback, or abandoning a session. Notably, such features are cumulatively impacted by the user’s long-term history, which makes RL especially challenging due to difficult credit assignment, where the impact of any individual action is usually small and noisy.¹

In this paper, we introduce *Dynamic Contextual Markov Decision Processes (DCMDPs)* to model such environment dynamics in a *history-dependent contextual* fashion. DCMDPs decompose the state space to include dynamic history-dependent contexts, where each context represents a different MDP, e.g., preferences of a human interacting with an agent, being affected by previous interactions. Particularly, we introduce a special class of *logistic DCMDPs*, in which context dynamics are determined by the aggregation of a set of feature vectors—functions of the immediate context, state and action—over time. This model is inspired by various psychological studies of human learning and conditioning; in particular, the Rescorla-Wagner (RW) model (Rescorla, 1972), a neuroscience model which describes the diminishing impact of repeated exposure to a stimulus due to historical conditioning. Critically, this structure allows us to develop tractable, UCB-style algorithms (Auer et al., 2008) for logistic DCMDPs that break the exponential dependence on history length in general DCMDPs.

Our contributions are as follows: (1) We introduce DCMDPs, a model that captures non-Markov context dynamics. (2) We introduce a subclass of DCMDPs for which state-action-context features are aggregated over time to determine context dynamics. We show how such problems can

^{*}Equal contribution ¹Google Research ²CREST, ENSAE. Correspondence to: Guy Tennenholtz <guytenn@gmail.com>.

¹A similar problem occurs in medical settings, where a patient’s previous reactions to certain treatments could implicitly affect the physician’s receptivity for treatment recommendations over long horizons. Another example includes human driver interventions in autonomous vehicles, where humans may take control of a vehicle for short periods of time.

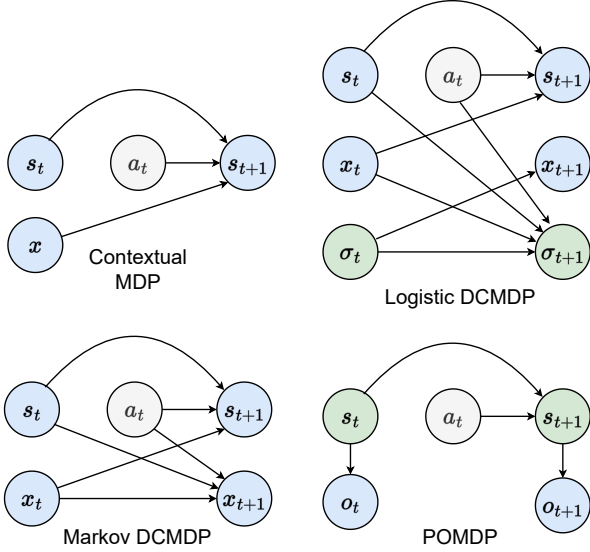


Figure 1: Causal diagrams comparing Contextual MDPs (Hallak et al., 2015), Markov DCMDPs (Section 2.1) and Logistic DCMDPs (Section 3). Logistic DCMDPs are history dependent, where $\sigma_t = \sum_{l=0}^{t-1} \alpha^{t-l-1} \mathbf{f}_l(s_l, a_l, x_l)$, and $\mathbf{f}_l : S \times A \times X \rightarrow \mathbb{R}^M$ are *unknown*, non-stationary, vector valued feature mappings. Green circles represent unobserved variables.

be solved by devising sample efficient and computationally tractable solutions, for which we establish regret bounds. (3) Inspired by our theoretical results, we construct a practical algorithm, based on MuZero (Schrittwieser et al., 2020), and demonstrate its effectiveness on a recommendation system benchmark with long history-dependent contexts.

2. Dynamic Contextual MDPs

We begin by defining *Dynamic Contextual MDPs (DCMDPs)*, a general framework for modeling history-dependent contexts². Let S, A and X be state, action, and context spaces, with cardinalities $|S|, |A|, |X|$, respectively. For any time $t \geq 1$, let $H_t = \mathcal{F}(s_1, a_1, x_1, \dots, s_t, a_t, x_t)$ be the set of histories up to time t ; and let $H = \bigcup_t H_t$. We denote $(s_0, a_0, x_0) = \cdot$.

A DCMDP is given by the tuple (X, S, A, r, P, H) , where, $r : S \times A \times X \rightarrow [0, 1]$ is a reward function, $P : H \times S \times A \rightarrow \Delta_S$ is a history-dependent transition function, and H is the horizon. DCMDP dynamics proceeds in discrete episodes $k = 1, 2, \dots, K$. At the beginning of episode k , the agent is initialized at state s_1^k . At any time h , the agent is in state s_h^k , has observed a history

²The term ‘‘context’’, as opposed to ‘‘state’’, differentiates between the Markov part of the state and the history dependent part of the state. Additionally, contexts often quantify characteristics of the environment (e.g., types of humans-in-the-loop), which can evolve in a distinct fashion, in contrast to the rest of the state.

$\tau_h^k = (s_1^k, a_1^k, x_1^k, \dots, s_{h-1}^k, a_{h-1}^k, x_{h-1}^k) \in H_h$, and selects an action $a_h^k \in A$. Then, the next context x_h^k occurs with (history-dependent) probability $P(x_h^k | \tau_h^k)$, the agent receives reward $r(s_h^k, a_h^k, x_h^k)$, and the environment transitions to state s_{h+1}^k with probability $P_h(s_{h+1}^k | s_h^k, a_h^k, x_h^k)$.

A policy $\pi : S \times H \rightarrow \Delta_A$ maps states and histories to distributions over actions. The value of π at time h is defined as $V_h^\pi(s, \tau) = \mathbb{E}_\pi \left[\sum_{t=h}^H r(s_t, a_t, x_t) \mid s_h = s, \tau_h = \tau \right]$, where $a_t = \pi(s_t, \tau_t)$, and $x_t = P(x_t | \tau_t)$. An optimal policy π^* maximizes the value over all states and histories; we denote its value function by V^* . We measure the performance of an RL agent by its *regret* – the difference between its value and that of an optimal policy: $\text{Reg}(K) = \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k)$.

Figure 1 depicts causal diagrams comparing general POMDPs to different types of DCMDPs, including three special cases: Contextual MDPs (Hallak et al., 2015), Markov DCMDPs, and logistic DCMDPs (defined in the next two sections). DCMDPs are closely related to POMDPs, yet their causal structure allows us to devise more tractable solution (characterized by an aggregation function, as we’ll see in Section 3) which can efficiently and tractably capture very long histories. In the next section, we describe a simple instance of DCMDPs, for which contexts are Markov, and show that standard MDP solutions can be applied. Then, in Section 3, we describe a more general DCMDP model, which uses aggregated features to represent histories, for which we provide sample efficient solutions and strong regret guarantees.

2.1. Markov DCMDPs

As a warm-up, we consider a simple version of DCMDPs in which context distributions are Markov w.r.t. the state and previous context. Specifically, we define a *Markov DCMDP* as a DCMDP which satisfies for all $h \geq [H]$, $\tau_h = (x_1, s_1, a_1, \dots, x_{h-1}, s_{h-1}, a_{h-1}) \in H_h$, $P(x_h | \tau_h) = P(x_h | s_{h-1}, a_{h-1}, x_{h-1})$. A Markov DCMDP $\mathcal{M} = (X, S, A, r, P, H)$ can be reduced to an MDP by augmenting the state space to include the context. To see this, we define the augmented MDP $\overline{\mathcal{M}} = (\overline{S}, A, \overline{r}, \overline{P}, H)$, where $\overline{S} = S \times X$ and $\overline{r}(\overline{s}_t, a_t) = r(s_t, a_t, x_t)$, $\overline{P}(\overline{s}_{t+1} | \overline{s}_t, a_t) = P(s_{t+1} | s_t, a_t, x_t) P(x_{t+1} | s_t, a_t, x_t)$. As a consequence, the Markov DCMDP \mathcal{M} and the MDP $\overline{\mathcal{M}}$ have the ‘‘same’’ optimal policy and value, and \mathcal{M} can be solved using standard RL methods, e.g., using UCBVI (Azar et al., 2017) one can obtain $\text{Reg}(K) = \tilde{O}\left(\sqrt{H^3 S A X K}\right)$. Markov DCMDPs also generalize contextual MDPs in an especially simple way; but they fail to capture the history dependence of contexts embodied by general DCMDPs. In the next section, we turn to a special case of DCMDPs that does so, but also admits tractable solution methods.

3. Logistic DCMDPs

We introduce a general class of DCMDPs, called *logistic DCMDPs*, where history dependence is structured using an aggregation of state-action-context-dependent features. Unlike Markov DCMDPs, logistic DCMDPs allow for context transitions to depend on history.

We define the softmax function $z_i : \mathbb{R}^M \rightarrow [0, 1]$, with temperature $\eta > 0$ as

$$z_i(\mathbf{u}) = \frac{\exp(\eta u_i)}{1 + \sum_{m=1}^M \exp(\eta u_m)} \quad (1)$$

for $i \in [M]$, $\mathbf{u} \in \mathbb{R}^M$, and $z_{M+1}(\mathbf{u}) = 1 - \sum_{i=1}^M z_i(\mathbf{u})$.

Definition 3.1 (Logistic DCMDP). A *logistic DCMDP* with latent feature maps $\{\mathbf{f}_h : S \times A \times \mathcal{X} \rightarrow \mathbb{R}^M\}_{h=0}^{H-1}$ is a DCMDP with context space $\mathcal{X} = \{x^{(i)}\}_{i=1}^{M+1}$, which satisfies, for all $h \in [H]$, $\tau_h = (s_1, a_1, x_1, \dots, s_{h-1}, a_{h-1}, x_{h-1}) \in \mathcal{H}_h$, and $i \in [M+1]$:

$$P_{\mathbf{f}}(x_h^{(i)} | \tau_h) = z_i \left(\sum_{t=0}^{h-1} \alpha^{h-t-1} \mathbf{f}_t(s_t, a_t, x_t) \right),$$

where $\alpha \in [0, 1]$ is a *history discount factor*.

Note that the latent functions \mathbf{f}_h are vector-valued and *unknown*. In a recommender system, \mathbf{f}_h may represent a user's unknown degree of trust in the system, or the effect of a sequence of recommendations on their satisfaction. The discount α allows for immediate effects to diminish over time (if less than 1).

A logistic DCMDP is denoted by $(X, S, A, r, P, H, \mathbf{f}, \alpha)$. We assume \mathbf{f} is ℓ_2 -bounded with $\sqrt{\sum f_{h,i}^2(s, a, x)} \leq L$, and we denote

$$F = \mathbf{f} \mathbf{f} : \prod_{h,i} f_{h,i}(s, a, x) \prod_{h,i} b_{h,i}(s, a, x) \quad (2)$$

the (rectangular) set where $b_{h,i}(s, a, x)$ are upper bounds on \mathbf{f} . Throughout our analysis we denote the effective history horizon $H_\alpha = \frac{\alpha^{2H}-1}{\alpha-1}$, and without loss of generality scale transitions in z_i (Equation (1)) with temperature $\eta = H_\alpha^{1/3}$. For clarity, we write $r(s, a, x^{(i)}) = r_i(s, a)$, $P(s^0 | s, a, x^{(i)}) = P_i(s^0 | s, a)$, and $\mathbf{r}(s, a) = (r_1(s, a), \dots, r_{M+1}(s, a))^T$, $\mathbf{P}(s^0 | s, a) = (P_1(s^0 | s, a), \dots, P_{M+1}(s^0 | s, a))^T$. We also denote by $n_h^k(s, a, x)$ the number of visits to s, a, x at time step h of episode $k-1$.

Next, we define a sufficient statistic for logistic DCMDPs that will prove valuable in our solution methods that follow.

³We set $\eta = H_\alpha^{1/2}$ for convenience. Different choices of η are equivalent to varying the bounds on F in Equation (2).

Definition 3.2 (Sufficient Statistic). Given a logistic DCMDP with feature maps \mathbf{f} , define $\sigma : \mathcal{H} \rightarrow \mathbb{R}^M$ as $\sigma(\tau_h; \mathbf{f}) := \sum_{t=0}^{h-1} \alpha^{h-t-1} \mathbf{f}_t(s_t, a_t, x_t)$, and the set of sufficient statistics by $\Sigma(\mathbf{f}) := \mathbf{f} \sigma(\tau; \mathbf{f}) \mathcal{G}_{\tau \in \mathcal{H}}$.

In Appendix B.1, we prove that $\sigma(\tau_h; \mathbf{f})$ is a sufficient statistic of the history for purposes of computing the optimal policy at time h . We do so by defining an equivalent MDP with state space $S = \Sigma(\mathbf{f})$ with well-defined dynamics and reward, and an equivalent optimal policy, which achieves the same optimal value.

Finally, similar to previous work on logistic and multinomial bandits (Abeille et al., 2021; Amani and Thrampoulidis, 2021), we define a problem-dependent constant for logistic DCMDPs which plays a key role in characterizing the behavior of $M+1$ multinomial logit bandit algorithms. For $\mathbf{x} \in \mathbb{R}^{M+1}$ and $\tau \in \mathcal{H}$, let $\mathbf{z}(\mathbf{x}) = (z_0(\mathbf{x}), \dots, z_{M+1}(\mathbf{x}))^T$, $\mathbf{A}(\tau; \mathbf{f}) = \text{diag}(\mathbf{z}(\sigma(\tau; \mathbf{f})))$, $\mathbf{z}(\sigma(\tau; \mathbf{f})) \mathbf{z}(\sigma(\tau; \mathbf{f}))^T$, and $1/\kappa = \inf_{\tau \in \mathcal{H}} \lambda_{\min} \mathbf{f} \mathbf{A}(\tau; \mathbf{f}) \mathbf{g}$. Informally, κ is related to saturation of the softmax z_i . For logistic DCMDPs, it is related to a worst-case context distribution w.r.t. \mathbf{f} and $\tau \in \mathcal{H}$. We refer to Abeille et al. (2021); Amani and Thrampoulidis (2021) for details, as well as lower bounds using this constant in logistic bandits.

The Rescorla-Wagner Model in Recommenders. Before continuing to provide sample efficient methods for solving logistic DCMDPs, we turn to motivate the aggregated model of history through the lens of the Rescorla-Wagner (RW) model (Rescorla, 1972) in a recommendation setting.

Logistic DCMDPs generate context transitions based on the sum of specific features of prior states, actions, and contexts, as captured by \mathbf{f} , with backward discounting to diminish the effect of past features or experiences, as captured by α . Such a model can be used to capture a (very simple) RW formulation of user behavior in an interactive recommender system. Let $I = \{i_1, \dots, i_n\}$ be a set of items. A user may like, dislike, or be unfamiliar with any of these items, represented by $u \in [-1, 0, 1]^n$. Let g_t be the user's (latent) current degree of satisfaction or engagement with the system. At each time t , the system asks the user for their disposition (e.g., rating) of an item $i_t \in I$. The user decides to answer the question with probability $z_1(g_t)$ (Equation (1)), which is strictly increasing with higher degrees of engagement level. The engagement level then evolves as $g_{t+1} = \alpha g_t + \beta u_{i_t}$, where $\alpha \in [0, 1]$, and β is a user-specific sensitivity factor. This model gives rise to a logistic DCMDP, whose solution gives the optimal recommender system policy. Specifically, actions $a_t := i_t \in I$ are the questions asked by the system, $f(s_t, a_t, x_t) = \beta u_{a_t}$ depends only on a_t , user engagement is $g_h = \sum_{t=0}^{h-1} \alpha^{h-t-1} f(s_t, a_t, x_t) = \sum_{t=0}^{h-1} \alpha^{h-t-1} \beta u_{i_t}$, x_t is the decision whether to answer, and s_t is the observation of the answer.

Algorithm 1 LDC-UCB

-
- 1: **for** $k = 1, \dots, K$ **do**
 - 2: $\hat{r}_{i,h}^k(s, a) = \hat{r}_{i,h}^k(s, a) + b_{i,h}^k(s, a), \delta i, h, s, a$
 - 3: $\hat{\pi}^k$ Optimistic Planner($\bar{M}_k(\delta)$) // Eq. 6
 - 4: Rollout a trajectory by acting $\hat{\pi}^k$
 - 5: $\hat{\mathbf{f}}^k \geq \arg \max_{\mathbf{f} \geq C_k(\delta)} L_{\lambda}^k(\mathbf{f})$ // Eq. 3
 - 6: Update $\hat{P}_i^{k+1}(s, a), \hat{r}_i^{k+1}(s, a), n^{k+1}(s, a, x)$ over rollout trajectory
 - 7: **end for**
-

4. Optimistic Methods for Logistic DCMDPs

Logistic DCMDPs' aggregation of features allow us to obtain sample efficient and computationally tractable solutions; namely, solutions which do not depend exponentially on history. In this section, we describe an optimistic algorithm for solving logistic DCMDPs and provide regret bounds. We focus on theoretical motivations here, and address computational tractability in the next section.

We first develop *Logistic Dynamic Context Upper Confidence Bound (LDC-UCB)*, a general RL method for logistic DCMDPs with unknown latent features (see Algorithm 1). At each episode k , LDC-UCB uses estimates of rewards

$$\hat{r}_{x,h}^k(s, a) = \frac{\sum_{k^0=1}^k \mathbb{1}\{x_h^{k^0}=x, s_h^{k^0}=s, a_h^{k^0}=a\} r_h^{k^0}}{n_h^k(s, a, x)}, \text{ transitions}$$

$$\hat{P}_{x,h}^k(s^0 | s, a) = \frac{\sum_{k^0=1}^k \mathbb{1}\{x_h^{k^0}=x, s_h^{k^0}=s, a_h^{k^0}=a, s_{h+1}^{k^0}=s^0\}}{n_h^k(s, a, x)}, \text{ and a}$$

projected estimate of $\hat{\mathbf{f}}$, calculated by maximizing the regularized log likelihood:

$$L_{\lambda}^k(\mathbf{f}) = \sum_{k^0=1}^k \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{1}\{x_h^k = i\} \ell_{i,h}^k(\mathbf{f}) - \lambda k \mathbf{f}^T \mathbf{K}_2 \mathbf{f}, \quad (3)$$

where $\ell_{i,h}^k(\mathbf{f}) = \log(z_i(\sigma(\tau_h^k; \mathbf{f})))$, $\lambda > 0$, and recall that $\sigma(\tau_h^k; \mathbf{f}) = \sum_{t=0}^h \alpha^{h-t} \mathbf{1}^T \mathbf{f}_t(s_t^k, a_t^k, x_t^k)$.

We account for uncertainty in these estimates by incorporating optimism. For rewards and transitions, we add a bonus term $b_{i,h}^k$ (see Appendix C for explicit definitions) to the estimated reward (line 2). To incorporate optimism in the latent features $\hat{\mathbf{f}}$, we build on results from multinomial logistic bandits (Amani and Thrampoulidis, 2021). Specifically, we derive a confidence bound over $\hat{\mathbf{f}}$, for which with probability at least $1 - \delta$

$$\left\| g_k(\mathbf{f}) - g_k(\hat{\mathbf{f}}_t) \right\|_{H_k^{-1}(\mathbf{f})} \leq \beta_k(\delta), \quad (4)$$

where $H_k(\mathbf{f}) = r_{\mathbf{f}}^2 L_{\lambda}^k(\mathbf{f})$, $g_k(\mathbf{f}) = r_{\mathbf{f}} L_{\lambda}^k(\mathbf{f}) + D_k$, $\beta_k(\delta) = \frac{M^{5/2} S A H}{\rho_{\lambda}} (\log(1 + \frac{k}{d\lambda}) + 2 \log(\frac{2}{\delta})) + \sqrt{\frac{\lambda}{4M}} + \frac{\lambda}{\rho_{\lambda} L}$. See Appendix G for exact expressions and a proof of the bound in Equation (4).

Next, we leverage the bound in Equation (4) to construct a feasible set of logistic DCMDPs. Specifically, we define the confidence set

$$C_k(\delta) = \left\{ \mathbf{f} \geq F : \left\| g_k(\mathbf{f}) - g_k(\hat{\mathbf{f}}_t) \right\|_{H_k^{-1}(\mathbf{f})} \leq \beta_k(\delta) \right\}. \quad (5)$$

and the following set of logistic DCMDPs:

$$\bar{M}_k(\delta) = \left\{ (X, S, A, \bar{r}, \hat{P}, H, \mathbf{f}, \alpha) : \mathbf{f} \geq C_k(\delta) \right\}. \quad (6)$$

The optimistic policy $\hat{\pi}^k$ (line 3) is that with greatest value over all DCMDPs in $\bar{M}_k(\delta)$, i.e., $\hat{\pi}^k$ corresponding to $\max_{\bar{m} \in \bar{M}_k(\delta)} V(s_1; \bar{m})$.

Combining the above, we prove the following regret guarantee for Algorithm 1.

Theorem 4.1. *Let $\lambda = \Theta(\frac{HM^{2.5}SA}{L})$. With probability at least $1 - \delta$, the regret of Algorithm 1 is*

$$\text{Reg}(K) \leq \tilde{O}\left(\frac{\rho}{H^6 M^{4.5} S^2 A^2 L^2 \kappa K}\right).$$

The proof of Theorem 4.1 can be found in Appendix C. We note that computing the optimistic policy over $\bar{M}_k(\delta)$ (line 3) is computationally difficult, especially due the history dependence of π on the accumulated latent features $\sum_{t=1}^h \alpha^{h-t} \mathbf{f}(s_t, a_t, x_t)$. We address this challenge next.

5. Mitigating Computational Complexity

In this section we show how to relax LDC-UCB (Algorithm 1) to mitigate its high computational complexity. Importantly, we maintain regret guarantees similar to those of Theorem 4.1 while obtaining an exponential improvement to computational cost. We later use these results to construct a practical model-based algorithm in Section 6.

To address the computational challenges of Algorithm 1, we focus on two problems. The first involves the set $C_k(\delta)$ (Equation (5) and line 5 of Algorithm 1) – where computation of the maximum likelihood constrained to set $C_k(\delta)$ is intractable. To address this, we prove that the constraint on the maximum likelihood estimator can be replaced by a simpler, rectangular set, enabling efficient calculation of the projected maximum likelihood. The second challenge is the complexity of the optimistic planner (Equation (6) and line 3 of Algorithm 1). To overcome this, we develop a *local* confidence bound, for every state-action-context triple (s, a, x) , and show it can be leveraged to design an optimistic planner, using a novel thresholding mechanism for optimism in logistic DCMDPs. Pseudocode for this tractable variant of LDC-UCB is presented in Algorithm 2.

5.1. A Tractable Estimator

We begin by constructing a tractable estimator for the latent feature maps \mathbf{f} which, instead of projecting to the set $C_k(\delta)$, solves for projected maximum likelihood on the rectangular set F (Equation (2)). Let $\gamma_k(\delta) = \left(2 + 2L \frac{\rho}{MH} + \sqrt{2(1+L)}\right)\beta_k + \sqrt{\frac{2(1+L)HM}{\lambda}}\beta_k^2(\delta)$. We define the tractable maximum likelihood estimator $\hat{\mathbf{f}}_T^k \in \arg \max_{\mathbf{f} \in F} L_\lambda^k(\mathbf{f})$, and have the following bound.

Lemma 5.1. *With probability at least $1 - \delta$, for all $k \geq [K]$,*

$$\left\| \hat{\mathbf{f}}_T^k - \mathbf{f} \right\|_{\mathbf{H}_k(\mathbf{f})} \leq \gamma_k(\delta). \quad (7)$$

The proof (see Appendix G.2) uses a convex relaxation of the set $C_k(\delta)$. Notice that the confidence region for $\hat{\mathbf{f}}_T^k$ is looser than that for $\hat{\mathbf{f}}^k$ (see Equation (4)), as $\beta_k(\delta) < \gamma_k(\delta)$. Nevertheless, its computation is tractable.

Next we can exploit the confidence bound in Equation (7) to construct a *local* bound for every state-action-context triple (s, a, x) using the number of visits to (s, a, x) , i.e., $n_h^k(s, a, x)$. The following result uses structural properties of logistic DCMDPs to achieve a local bound for $\hat{\mathbf{f}}_T^k$. Its proof generalizes the local confidence bound in Tennenholtz et al. (2022), and can be found in Appendix G.3.

Lemma 5.2 (Local Estimation Confidence Bound). *For any $\delta > 0$, with probability of at least $1 - \delta$, for all $k \geq [K]$, $h \geq [H]$, $i \geq [M]$ and $s, a, x \in S \times A \times X$, it holds that*

$$\left| (\hat{\mathbf{f}}_T^k(s, a, x))_{i,h} - (\mathbf{f}(s, a, x))_{i,h} \right| \leq \frac{2^{\rho_{\kappa}} \bar{\kappa} \gamma_k(\delta)}{\sqrt{n_h^k(s, a, x) + 4\lambda}}.$$

Lemma 5.2 allows one to reason about the unknown features locally for any visited (s, a, x) , a vital step toward an efficient optimistic planner. Indeed, as we see in the next section, the cost of planning in logistic DCMDPs can be reduced significantly using this bound.

5.2. Threshold Optimistic Planning

We now address the major computational challenge of Algorithm 1 – the complexity of optimistic planning (line 3 of Algorithm 1). To do this, we leverage the local bound in Lemma 5.2 and construct an optimistic planner using a novel threshold mechanism, as we describe next.

Recall the set of sufficient statistics $\Sigma(\mathbf{f}) = \bar{\mathbf{f}}\sigma(\tau; \mathbf{f})\mathcal{G}_{\tau \geq H}$ (Definition 3.2), which is a finite, vector-valued set with cardinality $|\Sigma(\mathbf{f})| = O((SAMH)^{MH})$, making planning in state space $S \times \Sigma$ exponentially hard. Consequently, searching for the optimistic DCMDP in the space of feature maps satisfying $\mathbf{f} \in C_k(\delta)$ (Equation (6)) requires searching over an exponentially large space.

We mitigate this problem exponentially, by leveraging the local confidence bound in Lemma 5.2. Let $B_k(\delta) \subseteq \mathbb{R}^M \times \mathbb{R}^M$ be the rectangular cuboid of all candidate confidence intervals satisfying the bound in Lemma 5.2. That is, $B_k(\delta)$ is the set of all M dimensional intervals $\left[\mathbf{l}_h^k(s, a, x), \mathbf{u}_h^k(s, a, x) \right]$, such that for all h, s, a, x , $\mathbf{f}_h(s, a, x) \in \left[\mathbf{l}_h^k(s, a, x), \mathbf{u}_h^k(s, a, x) \right]$, where $\mathbf{u}_h^k(s, a, x), \mathbf{l}_h^k(s, a, x) = \hat{\mathbf{f}}_T^k \left(\rho \frac{2^{\rho_{\kappa}} \bar{\kappa} \gamma_k(\delta)}{n_h^k(s, a, x^{(1)}) + 4\lambda}, \dots, \rho \frac{2^{\rho_{\kappa}} \bar{\kappa} \gamma_k(\delta)}{n_h^k(s, a, x^{(M)}) + 4\lambda} \right)^T$. In what follows, we identify key characteristics of the optimistic value when optimized over $B_k(\delta)$. Specifically, we show that an optimistic solution lies on the extreme points of $B_k(\delta)$, but more importantly, at one of M *specific extreme points*. This limits the search required by optimistic planning to a much smaller set, which can be approximated effectively in practice.

Optimism in intervals. Instead of augmenting the state space with $\Sigma(\mathbf{f})$, we use the set of confidence intervals defined by $B_k(\delta)$. We denote by $\mathcal{C}_h^k : \Sigma(\hat{\mathbf{f}}^k) \rightarrow \mathbb{R}^M \times \mathbb{R}^M$ the confidence interval of the sufficient statistic $\sigma(\tau_h^k, \hat{\mathbf{f}}^k)$. That is,

$$\begin{aligned} \mathcal{C}_h^k &= \mathcal{C}(\sigma(\tau_h^k; \hat{\mathbf{f}}^k)) \\ &= \left[\sum_{t=0}^{h-1} \alpha^{h-t-1} \mathbf{l}_h^k(s_t^k, a_t^k, x_t^k), \sum_{t=0}^{h-1} \alpha^{h-t-1} \mathbf{u}_h^k(s_t^k, a_t^k, x_t^k) \right]. \end{aligned}$$

We also denote by $\mathcal{I}^k = \left\{ \mathcal{C}(\sigma(\tau, \hat{\mathbf{f}}_T^k)) \right\}_{\tau \geq H}$ the set of possible confidence intervals over $B_k(\delta)$ in episode k .

Next, we augment the state space S at every episode k by $S \times \mathcal{I}^k$, and define the augmented state-action optimistic value for context $i \geq [M+1]$ and confidence interval $\mathcal{C}_h^k = \mathcal{C}(\sigma(\tau_h^k, \hat{\mathbf{f}}^k))$ at time step $h \geq [H]$ by

$$\bar{Q}_i(s, a, \mathcal{C}_h^k) = \bar{r}_i(s, a) + \mathbb{E}_{s^0} \mathbb{P}_i(j_{s,a}) \left[\bar{V}_{h+1}(s^0, \mathcal{C}_{h+1}^k) \right],$$

where, with slight abuse of notation, we used $\mathcal{C}_{h+1}^k = \mathcal{C}(\sigma(\tau_h^k \llbracket \{s, a, x^{(i)}\}, \hat{\mathbf{f}}_T^k \rrbracket))$ to denote the next aggregated confidence interval. The optimistic value \bar{V}_h is defined by maximizing over sufficient statistics in the confidence set \mathcal{C}_h^k and $a \in A$. That is,

$$\bar{V}_h(s, \mathcal{C}_h^k) = \max_{a \in A} \max_{\bar{\sigma} \in \mathcal{C}_h^k} \sum_{i=1}^{M+1} z_i(\bar{\sigma}) Q_i(s, a, \mathcal{C}_h^k) \quad (8)$$

Indeed, \bar{V}_h is an optimistic value, as shown by the following proposition. Its proof is provided in Appendix D.3.

Proposition 5.3 (Optimistic Value). *Let \bar{V}_h as defined in Equation (8). Then, w.h.p. $\bar{V}_1(s_1^k, \mathcal{C}_1^k) \geq V_1(s_1^k)$.*

Algorithm 2 Tractable LDC-UCB

-
- 1: **for** $k = 1, \dots, K$ **do**
 - 2: $\hat{r}_{i,h}^k(s, a) = \hat{r}_{i,h}^k(s, a) + b_{i,h}^k(s, a), \delta_i, h, s, a$
 - 3: $\hat{\pi}^k$ Optimistic DP($\hat{r}^k, \hat{P}^k, B_k(\delta)$) // Eq. 8
 - 4: Rollout a trajectory by acting $\hat{\pi}^k$
 - 5: $\hat{\mathbf{f}}_T^{k+1} \geq \arg \max_{\mathbf{f} \geq \mathbf{F}} L_\lambda^k(\mathbf{f})$ // Eq. 3
 - 6: Update $\hat{P}_i^{k+1}(s, a), \hat{r}_i^{k+1}(s, a), n^{k+1}(s, a, x)$ over rollout trajectory
 - 7: **end for**
-

Next, we turn to show that the maximization problem in Equation (8) can be solved efficiently, though \mathbf{C}_h^k is an exponentially large set. Notice that the inner term $\sum_{i=0}^M z_i(\bar{\sigma}) Q_i(s, a, \mathbf{C}_h^k)$ in Equation (8) is not convex. Still, our analysis shows that a solution to the inner maximization problem lies in the set of extreme points of \mathbf{C}_h^k . That said, these 2^M extreme points make exhaustive search intractable. Fortunately, we can also show that the optimal solution lies in a space of exactly M solutions – a linearly sized, tractable search space.

To this end, we define the threshold set, which we will use to construct the (linear) set of feasible extreme points.

Definition 5.4. For a rectangular cuboid defined by the interval $\mathbf{C} = [\mathbf{l}, \mathbf{u}] \in \mathbb{R}^{M+1} \times \mathbb{R}^{M+1}$, vector $\mathbf{y} \in \mathbb{R}^{M+1}$ and real number $t \in \mathbb{R}$ we define $\mathbf{th}_t(\mathbf{y}, \mathbf{C}) \in \mathbb{R}^{M+1}$ by

$$[\mathbf{th}_t(\mathbf{y}, \mathbf{C})]_i = \begin{cases} l_i & y_i < t \\ u_i & \text{o.w.} \end{cases}$$

Definition 5.5. For a vector $\mathbf{Q} \in \mathbb{R}^{M+1}$, we define the threshold set $T(\mathbf{Q}) = \left\{ \frac{Q_i + Q_{i+1}}{2} \right\}_{i=1}^M$.

We use these definitions to show that the optimal solution to Equation (8) lies in the threshold set of Q -values (see proof in Appendix F.1).

Lemma 5.6 (Threshold Optimism). *Let $\mathbf{Q} \in \mathbb{R}^{M+1}$. For any $\mathbf{x} \in \mathbb{R}^{M+1}$ such that $x_i = 0$ define $f(\mathbf{x}) = \sum_{i=1}^{M+1} z_i(\mathbf{x}) Q_i$. Let $\mathbf{C} = [\mathbf{l}, \mathbf{u}] \in \mathbb{R}^{M+1} \times \mathbb{R}^{M+1}$ and assume that $\mathbf{l} < \mathbf{u}$. Then, there exists $t \in T(\mathbf{Q})$ such that $\mathbf{th}_t(\mathbf{Q}, \mathbf{C}) \geq \arg \max_{\mathbf{x} \in \mathbf{C}} f(\mathbf{x})$.*

We can now leverage Lemma 5.6 to solve the inner maximization in Equation (8). For notational convenience, we write $\bar{Q}_i = \bar{Q}_i(s, a, \mathbf{C}_h^k)$ and $\mathbf{Q} = (Q_1, \dots, Q_{M+1})^T$. Applying Lemma 5.6, we get that

$$\max_{\bar{\sigma} \in \mathbf{C}_h^k} \sum_{i=1}^{M+1} z_i(\bar{\sigma}) \bar{Q}_i = \max_{t \in T(\mathbf{Q})} \sum_{i=1}^{M+1} z_i(\mathbf{th}_t(\mathbf{Q}, \mathbf{C}_h^k)) \bar{Q}_i. \quad (9)$$

As a result, the non-convex maximization problem in Equation (9) reduces the search space to M optimistic candidates.

5.3. Putting It All Together

Using Lemma 5.6 and particularly its derived corollary in Equation (9), we construct an optimistic planner, denoted by Optimistic DP, which plans via dynamic programming using Equation (9); we refer to Appendix F for an explicit formulation of the optimistic planner. Finally, using the tractable estimator $\hat{\mathbf{f}}_T^k$, and the threshold optimistic planner, we present a tractable variant of LDC-UCB in Algorithm 2, for which we have the following regret guarantee.

Theorem 5.7. *Let $\lambda = \Theta\left(\frac{HM^{2.5}SA}{L}\right)$. With probability at least $1 - \delta$, the regret of Algorithm 2 is*

$$R(K) \leq \tilde{O}\left(\frac{D}{H^8 M^{6.5} S^2 A^2 L^4 \kappa K}\right).$$

The proof of the theorem can be found in Appendix D. As expected, the regret upper bound in Algorithm 2 is worse than that of Algorithm 1 by a factor of $\tilde{O}(HML)$. This result is strongly affected by the looser bound for the tractable feature maps in Lemma 5.2. Nevertheless, the intractability of Algorithm 1 compared to the tractability of Algorithm 2 suggests this is a more-than-reasonable tradeoff. Moreover, our tractable variant of LDC-UCB gives rise to practical optimistic algorithms, as we demonstrate next.

6. DCZero

Motivated by our theoretical results, we present a practical model-based optimistic algorithm for solving DCMDPs. We build on MuZero (Schrittwieser et al., 2020), a recent model-based algorithm which constructs a model in latent space and acts using Monte Carlo Tree Search (MCTS, Coulom (2007)). MuZero uses representation, transition, and prediction networks for training and acting. The representation network first embeds observations in a latent space, after which planning takes place using the transition and prediction networks through a variant of MCTS. Importantly, instead of predicting the next state (e.g., using world models (Hafner et al., 2023)), MuZero trains its latent space by predicting three quantities—the reward, value, and current policy—by rolling out trajectories in latent space (see Schrittwieser et al. (2020) for further details).

We develop DCZero, an algorithm based on MuZero for DCMDPs (see Algorithm 3). Like MuZero, DCZero uses representation, transition, and prediction networks to learn and act in the environment. In contrast to MuZero, DCZero trains an additional ensemble of networks to estimate the unknown features \mathbf{f} using cross-entropy. Estimated quantities of the ensemble are used to construct confidence intervals for the sufficient statistics, which are used to augment the state. DCZero uses $M + 1$ transition networks (one for each context), and predicts $M + 1$ reward functions. To incorporate optimism, the value function is trained optimistically using the thresholding technique in the previous

section, where rewards for unseen actions are sampled from the trained reward models r_i and next states are sampled from the trained models P_i .

Movie Recommendation Environment. To evaluate the effectiveness of DCZero, we develop a movie recommendation environment based on the MovieLens dataset (Harper and Konstan, 2015). Users and items are represented in embedding space computed using SVD of the MovieLens ratings matrix. Each of n users is assigned a set of M possible user embeddings; i.e., each user $u \in \{u^{(i)}\}_{i=1}^n$ is assigned a set of preference vectors $\mathbf{x} = \{\mathbf{x}^{(j)}\}_{j=1}^{M+1}$, $\mathbf{x}^{(j)} \in \mathbb{R}^d$. Intuitively, these vectors reflect distinct user preferences corresponding to some aspect of the user’s latent state (e.g., mood or current interest (Cen et al., 2020); location, companions, or activity; level of trust or satisfaction with the system) and hence influence u ’s behavior.

The recommendation agent interacting with a user selects an item x from a random set of A movies, $\hat{v}^{(a)} \mathcal{G}_{a=1}^A$, $\mathbf{v}^{(a)} \in \mathbb{R}^d$, and recommends it. The user context then evolves according to some history-dependent dynamics represented by a logistic DCMDP. Specifically, we assume unknown latent features $\mathbf{f}(\mathbf{x}, \mathbf{v})$ with the user’s aggregated features (at time $h \in [H]$, episode k) being: $\boldsymbol{\sigma}_{k,h} = \sum_{t=0}^{h-1} \alpha^{h-t-1} \mathbf{f}(\mathbf{x}_k^{(j_t)}, \mathbf{v}^{(a_t)})$. The agent recommends movie $\mathbf{v}^{(a)}$ to the user, while the user preference vector is sampled as $\mathbf{x}_k^{j_h} \sim z(\boldsymbol{\sigma}_{k,h})$. The agent then receives a reward $r_j(\mathbf{x}, a) = (\mathbf{x}_k^{(j_h)})^T \boldsymbol{\Sigma} \mathbf{v}^{(a)}$ reflecting the user’s (current) preference for the movie, and the user’s latent state transitions given the unknown function $\mathbf{f}(\mathbf{x}, \mathbf{v})$ and discount α ; that is, $\boldsymbol{\sigma}_{k,h+1} = \alpha \boldsymbol{\sigma}_{k,h} + \mathbf{f}(\mathbf{x}_k^{(j_h)}, \mathbf{v}^{(a)})$.

We test our methods in two variants of this environment. In the first, “AttractionEnv”, user latent features \mathbf{f} are correlated with the user’s degree of preference for the recommended movie:

$$\mathbf{f}(\mathbf{x}^{(j)}, \mathbf{v}) = \mu((\mathbf{x}^{(j)})^T \boldsymbol{\Sigma} \mathbf{v}), \quad (\text{Attraction})$$

where μ is a component-wise monotonically increasing function. AttractionEnv reflects users with a tendency to desire content similar to those they most recently consumed. This may reflect the positive influence of exposure to new types of content, increased familiarity increasing preference, or content domains (such as music) where some mild consistency of experience is preferred to jarring shifts in style or genre. The second environment, “NoveltyEnv”, reflects a contrasting dynamics in which user latent features evolve such that \mathbf{f} is anti-correlated with the user’s preference for the recommended movie:

$$\mathbf{f}_i(\mathbf{x}^{(j)}, \mathbf{v}) = \begin{cases} \mu((\mathbf{x}^{(j)})^T \boldsymbol{\Sigma} \mathbf{v}) & , j = i \\ \mu((\mathbf{x}^{(j)})^T \boldsymbol{\Sigma} \mathbf{v}) & , \text{o.w.} \end{cases} \quad (\text{Novelty})$$

Algorithm 3 DCZero

- 1: **require:** Size of ensemble B
 - 2: **init:** Replay buffer \mathcal{R} ;
 - 3: **for** $k = 1, 2, \dots$ **do**
 - 4: Train bootstrap ensemble of B feature maps $\left\{ \hat{\mathbf{f}}_{\theta_b} : S \times A \times \mathcal{V} \rightarrow \mathbb{R}^M \right\}_{b=1}^B$ over \mathcal{R} using cross-entropy loss.
 - 5: Augment state s with aggregated feature confidence $\mathbf{C}_h = \sum_{t=0}^{h-1} \alpha^{h-t-1} \text{std}_b \left(\left\{ \hat{\mathbf{f}}_{\theta_b}(s_t, a_t, x_t) \right\}_{b=1}^B \right)$.
 - 6: Train threshold optimistic value estimator over \mathcal{R}
$$Q_{i,\psi}(s, a, \mathbf{C}_h) = \hat{r}_i(s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}_i(s,a)} V_\psi(s', \mathbf{C}_{h+1}),$$

$$V_\psi(s, \mathbf{C}_h) = \max_{t \in T(\mathbf{Q}_\psi)} \sum_{i=1}^{M+1} z_i \left(\mathbf{th}_t(\mathbf{Q}_\psi, \mathbf{C}_h) \right) Q_{i,\psi}^{\pi_\phi(s)}.$$
 - 7: Act and train representation network, M transition networks (for each \hat{P}_i), and M prediction networks (for each \hat{r}_i) using MuZero-ALG with the optimistic value \bar{V}_ψ in MCTS. Return replay buffer \mathcal{R} .
 - 8: **end for**
-

As a result, movies that previously appealed to the user become less preferred, reflecting a desire for novelty over short time periods.

Experiments. All experiments used a horizon of $H = 300$, $M = 6$ user classes, $A = 6$ slate items (changing every reset), and a user embedding dimension of $d = 20$. We used default parameters for MuZero and applied the same parameters to DCZero. We compared DCZero and MuZero on the AttractionEnv and NoveltyEnv environments. We also tested a history-dependent variant of MuZero, which uses the sequence of past movies and contexts to densely represent history. More specifically, Hist-MuZero uses a stack of 30 previous observations as its state. We implemented both MLP and Transformer-based model architectures, but present results for the Transformer, as both had similar performance.

Figure 2 shows these comparisons. The plots compare the return of DCZero with the two baselines on AttractionEnv and NoveltyEnv with $\alpha = 0.99$; we also vary the values of α on the AttractionEnv. We see that DCZero is able to outperform both baselines, with significant increases in performance for larger values of α (i.e., longer history dependence). This suggests that DCZero can be especially beneficial in problems that exhibit long history dependence. Interestingly, we note that using a dense history-dependent Transformer hurts performance, except for very small values of α (indeed, only for $\alpha = 0.1$ does the sequence model outperform the other methods).

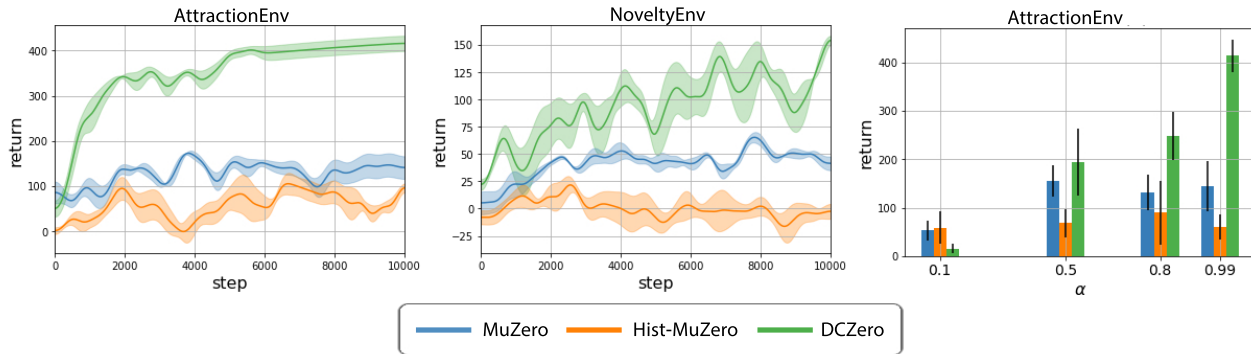


Figure 2: Plots comparing MuZero, Hist-MuZero, and DCZero on the AttractionEnv(left) and NoveltyEnv (middle). We also compare results for different values of α (right). All experiments show mean scores with 95% confidence intervals.

7. Related Work

Contextual MDPs (Hallak et al., 2015) have proven useful in a numerous studies (Jiang et al., 2017; Zintgraf et al., 2019; Kwon et al., 2021). Contexts are sampled once and are fixed throughout the episode. DCMDPs can be seen as a generalization of contextual MDPs, where contexts can change over time in a realistic, history-dependent fashion. Other forms of DCMDPs, are interesting directions for future work, including DCMDPs for which contexts change slowly in time. In Chen et al. (2022) a latent context variable changes abruptly at discrete points in time. Our logistic DCMDP considers history-dependent dynamics of contexts, which can depend on previous states and actions. Moreover, our model can capture smoother behavior which changes very slowly over time (over long histories). Finally, in contrast to Chen et al. (2022), our work provides theoretical guarantees, showing statistical and computational efficiency of our approach. In Mao et al. (2018), a non-stationary contextual environment is considered, yet the contexts are not allowed to depend on previous states and actions. Ren et al. (2022) propose a Bayesian approach for learning contextual MDPs for which contexts can change dynamically. Nevertheless, their model assumes dynamics that are not state-action dependent, and not history dependent.

Partially observable MDPs are widely studied (Papadimitriou and Tsitsiklis, 1987; Vlassis et al., 2012; Krishnamurthy et al., 2016; Tennenholtz et al., 2020; Xiong et al., 2022). As POMDPs are inherently history dependent, recent work has identified models and assumptions for which sample-efficient algorithms can be derived (Xiong et al., 2022; Liu et al., 2022a;b). Nevertheless, such solutions are often computationally intractable, impeding their practical implementation. With DCMPDs, we focus on specific forms of history-dependence, and show them to be computationally tractable, as well as effectively deployable.

Tennenholtz et al. (2022) define TerMDPs, a framework which models exogenous, non-Markovian termination in the

environment. Once terminated, the agent stops acting and accrues no further rewards. TerMDPs capture various scenarios in which exogenous actors disengage with the agent (e.g., passengers in autonomous vehicles or users abandoning a recommender), and can be shown to be a special case of logistic DCMDPs (see Appendix B.2). As such, logistic DCMDPs support reasoning about optimizing more general contextual behavior, including: those involving notions of trust (e.g., where users become more or less receptive to agent recommendations); situations where humans override an agent for short periods; and modeling the effects of user satisfaction, moods, etc.

8. Discussion and Future Work

In this work we presented DCMDPs, and logistic DCMDPs in particular—a general history-dependent contextual framework which admits sample and computationally efficient solutions. The aggregation structure of logistic DCMDPs gives rise to efficient estimation of the unknown feature maps. We provided regret guarantees and developed a tractable realization of LDC-UCB using a computational estimator and a novel planning procedure. Finally, we tested DCZero, a model-based implementation of LDC-UCB, demonstrating its efficacy on a recommendation benchmark.

While logistic DCMDPs assume linear aggregations of past features, other variants with more complex parametric function classes over history are possible. Nevertheless, such complex function classes often require sample-inefficient techniques, suggesting that logistic DCMDPs may be especially well-suited to capturing extended, long history dependence. In particular, they admit sample and computationally efficient solutions, which can be implemented in practice. As future work, a hybrid approach which considers combining dense models (such as Transformers) for short-history dependence, and aggregated models (such as logistic DCMDPs) for very long history dependence, may offer the “best of both worlds” in practice.

Acknowledgements

This project has received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 101034255.



Nadav Merlis is partially supported by the Viterbi Fellowship, Technion.

References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- Marc Abeille, Louis Fauray, and Clément Calauzènes. Instance-wise minimax-optimal algorithms for logistic bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 3691–3699. PMLR, 2021.
- M Mehdi Afsar, Trafford Crump, and Behrouz Far. Reinforcement learning based recommender systems: A survey. *ACM Computing Surveys*, 55(7):1–38, 2022.
- Sanae Amani and Christos Thrampoulidis. Ucb-based algorithms for multinomial logistic regression bandits. *Advances in Neural Information Processing Systems*, 34: 2913–2924, 2021.
- Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement learning. *Advances in neural information processing systems*, 21, 2008.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR, 2017.
- Fahiem Bacchus, Craig Boutilier, and Adam Grove. Rewarding behaviors. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence (AAAI-96)*, pages 1160–1167, Portland, OR, 1996.
- KV Bhagwat and R Subramanian. Inequalities between means of positive operators. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 83, pages 393–401. Cambridge University Press, 1978.
- Yukuo Cen, Jianwei Zhang, Xu Zou, Chang Zhou, Hongxia Yang, and Jie Tang. Controllable multi-interest framework for recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD-20)*, pages 2942–2951, 2020.
- Li Chen and Pearl Pu. Critiquing-based recommenders: Survey and emerging trends. *User Modeling and User-Adapted Interaction*, 22(1):125–150, 2012.
- Xiaoyu Chen, Xiangming Zhu, Yufeng Zheng, Pushi Zhang, Li Zhao, Wenxue Cheng, Peng CHENG, Yongqiang Xiong, Tao Qin, Jianyu Chen, et al. An adaptive deep rl method for non-stationary environments with piecewise stable context. In *Advances in Neural Information Processing Systems*, 2022.
- Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International conference on computers and games*, pages 72–83. Springer, 2007.
- Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 5713–5723, 2017.
- Yonathan Efroni, Nadav Merlis, and Shie Mannor. Reinforcement learning with trajectory feedback. *arXiv preprint arXiv:2008.06036*, 2020.
- Yonathan Efroni, Nadav Merlis, Aadirupa Saha, and Shie Mannor. Confidence-budget matching for sequential budgeted learning. In *International Conference on Machine Learning*, pages 2937–2947. PMLR, 2021.
- Abdur R Fayjie, Sabir Hossain, Doukhi Oualid, and Deok-Jin Lee. Driverless car: Autonomous driving using deep reinforcement learning in urban environment. In *2018 15th international conference on ubiquitous robots (ur)*, pages 896–901. IEEE, 2018.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes. *arXiv preprint arXiv:1502.02259*, 2015.
- F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- Henning Hohnhold, Deirdre O’Brien, and Diane Tang. Focusing on the long-term: It’s good for users and business. In *Proceedings of the Twenty-first ACM International Conference on Knowledge Discovery and Data Mining (KDD-15)*, pages 1849–1858, Sydney, 2015.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low bellman rank are pac-learnable. In *International Conference on Machine Learning*, pages 1704–1713. PMLR, 2017.

- Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Pac reinforcement learning with rich observations. *Advances in Neural Information Processing Systems*, 29, 2016.
- Jeongyeol Kwon, Yonathan Efroni, Constantine Caramanis, and Shie Mannor. RL for latent mdps: Regret guarantees and a lower bound. *Advances in Neural Information Processing Systems*, 34:24523–24534, 2021.
- Qinghua Liu, Alan Chung, Csaba Szepesvári, and Chi Jin. When is partially observable reinforcement learning not scary? *arXiv preprint arXiv:2204.08967*, 2022a.
- Qinghua Liu, Praneeth Netrapalli, Csaba Szepesvari, and Chi Jin. Optimistic mle—a generic model-based algorithm for partially observable sequential decision making. *arXiv preprint arXiv:2209.14997*, 2022b.
- Hongzi Mao, Shaileshh Bojja Venkatakrishnan, Malte Schwarzkopf, and Mohammad Alizadeh. Variance reduction for reinforcement learning in input-driven environments. In *International Conference on Learning Representations*, 2018.
- Martin Mladenov, Ofer Meshi, Jayden Ooi, Dale Schuurmans, and Craig Boutilier. Advantage amplification in slowly evolving latent-state environments. In *Proceedings of the Twenty-eighth International Joint Conference on Artificial Intelligence (IJCAI-19)*, pages 3165–3172, Macau, 2019.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Christos H Papadimitriou and John N Tsitsiklis. The complexity of markov decision processes. *Mathematics of operations research*, 12(3):441–450, 1987.
- Hang Ren, Aivar Sootla, Taher Jafferjee, Junxiao Shen, Jun Wang, and Haitham Bou Ammar. Reinforcement learning in presence of discrete markovian context evolution. In *International Conference on Learning Representations*, 2022.
- Robert A Rescorla. A theory of pavlovian conditioning: Variations in the effectiveness of reinforcement and non-reinforcement. *Current research and theory*, pages 64–99, 1972.
- Alessandro Ronca and Giuseppe De Giacomo. Efficient pac reinforcement learning in regular decision processes. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI-21)*, pages 2026–2032, Montreal, 2021.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- Tianxiao Sun and Quoc Tran-Dinh. Generalized self-concordant functions: a recipe for newton-type methods. *Mathematical Programming*, 178(1):145–213, 2019.
- Guy Tennenholtz, Uri Shalit, and Shie Mannor. Off-policy evaluation in partially observable environments. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10276–10283, 2020.
- Guy Tennenholtz, Nadav Merlis, Lior Shani, Shie Mannor, Uri Shalit, Gal Chechik, Assaf Hallak, and Gal Dalal. Reinforcement learning with a terminator. In *Advances in Neural Information Processing Systems*, 2022.
- Chen Tessler, Guy Tennenholtz, and Shie Mannor. Distributional policy optimization: An alternative approach for continuous control. *Advances in Neural Information Processing Systems*, 32, 2019.
- Nikos Vlassis, Michael L Littman, and David Barber. On the computational complexity of stochastic controller optimization in pomdps. *ACM Transactions on Computation Theory (TOCT)*, 4(4):1–8, 2012.
- Tsachy Weissman, Erik Ordentlich, Gadiel Seroussi, Sergio Verdu, and Marcelo J Weinberger. Inequalities for the l_1 deviation of the empirical distribution. *Hewlett-Packard Labs, Tech. Rep*, 2003.
- Mark Wilhelm, Ajith Ramanathan, Alexander Bonomo, Sagar Jain, Ed H. Chi, and Jennifer Gillenwater. Practical diversified recommendations on YouTube with determinantal point processes. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM18)*, pages 2165–2173, Torino, Italy, 2018.
- Yi Xiong, Ningyuan Chen, Xuefeng Gao, and Xiang Zhou. Sublinear regret for learning pomdps. *Production and Operations Management*, 31(9):3491–3504, 2022.
- Xiaoxue Zhao, Weinan Zhang, and Jun Wang. Interactive collaborative filtering. In *Proceedings of the 22nd ACM International Conference on Information and Knowledge Management*, pages 1411–1420, 2013.
- Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. *arXiv preprint arXiv:1910.08348*, 2019.

Appendix

The appendix is organized as follows. First, in Appendix A, we define additional notations that are used throughout this work. We then show that σ is indeed a sufficient statistic for calculating the optimal policy in a logistic DCMDP (Appendix B.1). Next, we provide regret guarantees for our algorithms for solving logistic DCMDPs: in Appendices C and D, we bound the regret of LDC-UCB (Theorem 4.1) and its tractable variant (Theorem 5.7), respectively. Finally, Appendices E to G contain technical lemmas which are crucial for deriving the above regret guarantees. Specifically, Appendix E is dedicated for optimism lemmas and decomposing the regret for logistic DCMDPs; Appendix F deals with the threshold optimistic planning required for the tractable version of LDC-UCB; and Appendix G provides confidence sets for the regularized log likelihood procedure, following Abeille et al. (2021); Amani and Thrampoulidis (2021).

A. Additional Notations

In this section, we define additional notation that will be of use throughout the proof. We work with the natural filtration

$$F_k = \sigma\left(\left\{\tau_{H+1}^{k^0}\right\}_{k^0 \geq [k]}, s_1^{k+1}\right) = \sigma\left(\{(s_h^1, a_h^1, x_h^1, R_h^1)\}_{h=1}^H, \dots, \{(s_h^k, a_h^k, x_h^1, R_h^k)\}_{h=1}^H, s_1^{k+1}\right),$$

and notice that the policy π^k , which might depend on s_1^k , is F_{k-1} -measurable. For brevity, for any episode $k \geq [K]$ and time step $h \geq [H]$, we define the probability distribution over the contexts by $z_h^k = \mathbf{z}(\sigma(\tau_h^k; \mathbf{f}))$, with $z_{i,h}^k = z_i(\sigma(\tau_h^k; \mathbf{f}))$ for any $i \geq X$.

With slight abuse of notation, we treat the latent features \mathbf{f} as vectors in $F = \mathbb{R}^{S \times A \times [H] \times [M]}$ instead of a mapping $\mathbf{f}_h : S \times A \times X \times \mathcal{V} \times \mathbb{R}^M, \delta h \geq [H]$ and use the notations $\mathbf{f}_i(s, a, h, x) = \mathbf{f}(s, a, x, h, i)$. We also let \mathbf{d}_h^k be the empirical discounted visitation vector at episode k up to time step h , i.e.,

$$\mathbf{d}_h^k(s, a, t, x) = H_\alpha^{-1/2} \alpha^{h-t-1} \mathbb{1}\{s_t^k = s, a_t^k = a, x_t^k = x\},$$

where $H_\alpha = \frac{1 - \alpha^{2H}}{1 - \alpha} = \min\left\{H, \frac{1}{1 - \alpha}\right\}$ is a normalization factor describing the effective historical horizon. Then, one can write $\sigma_i(\tau_h^k; \mathbf{f}) = h \mathbf{f}_i, \mathbf{d}_h^k(s, a, t, x) i$. Notice that \mathbf{d}_h^k is a vector containing zeros except for h elements with the values $f H_\alpha^{-1/2}, H_\alpha^{-1/2} \alpha, \dots, H_\alpha^{-1/2} \alpha^{h-1} g$, where each value appears exactly once. We denote the set of all possible vectors of such form for any $h \geq [H]$ by \mathcal{D} , and notice that for all $\mathbf{d} \in \mathcal{D}$,

$$\|\mathbf{d}\|_2^2 = \sum_{h=1}^H H_\alpha^{-1} \alpha^{2h} = \frac{1 - \alpha^{2H}}{1 - \alpha} H_\alpha^{-1} = 1.$$

Next, we define the following summation operators:

For any fixed $h \geq [H]$ and $i \geq X$, if $P_{i,h} : S \times A \times \mathcal{V} \times \Delta_S$ is a transition kernel and $V : S \times H_h \times \mathcal{V} \times \mathbb{R}$ is a value function, the expected value is denoted by

$$[P_{i,h} V](s, a, \tau_h) = P_{i,h}(j_s, a)^T V(\cdot, \tau_h) = \sum_{s^0 \in \mathcal{S}} P_i(s_{h+1} = s^0 | s_h = s, a_h = a) V(s^0, \tau_h).$$

and, in general, use similar notations for any transition kernel $P : \mathcal{Y} \times \mathcal{V} \times \Delta_S$ from arbitrary space \mathcal{Y} .

We denote the vectorized version of $P_{i,h} V$ by

$$P_h V = (P_{1,h} V, \dots, P_{M,h} V)^T.$$

If $Z : \mathcal{Y} \times \mathcal{V} \times \Delta_X$ is a mapping to the probability simplex over X and $U : Z \times X \times \mathcal{V} \times \mathbb{R}$, where \mathcal{Y}, \mathcal{W} are some arbitrary spaces, we let

$$[ZU](y, w) = \sum_{i=1}^{M+1} Z(y, i) U(w, i),$$

and in particular, we use $Z_h^{\mathbf{f}}(\tau_h) = \mathbf{z}(\sigma(\tau_h; \mathbf{f}))$ and $Z_h^k = \mathbf{z}_h^k$.

Finally, given a transition kernel P and latent feature \mathbf{f} , we denote the transition operator over a value V by

$$[T_h^{P, \mathbf{f}} V](s, a, \tau_h) = [Z_h^{\mathbf{f}} P_h V](s, a, \tau_h) = \sum_{s^0 \in \mathcal{S}} \sum_{i=1}^M P_i(s^0 | s, a) z_i(\sigma(\tau_h; \mathbf{f})) V(s^0, \tau_h).$$

we similarly use the notation $T_h^{P, Z}$ when for general context distribution that are not necessarily by latent features \mathbf{f} .

B. Logistic DCMDPs

B.1. Sufficient Statistic

We prove that the σ is a sufficient statistic for calculating the optimal policy. We begin by defining an augmented MDP $\mathcal{M}_{\text{aug}} = (S_{\text{aug}}, A_{\text{aug}}, P_{\text{aug}}, r_{\text{aug}}, H)$, where $S_{\text{aug}} = S \times \Sigma(\mathbf{f})$ is the augmented state space, and $A_{\text{aug}} = A$ is the (unchanged) action space. The augmented transition function is defined for $s, \sigma \in S \times \Sigma(\mathbf{f}), a \in A, s^0, \sigma^0 \in S \times \Sigma(\mathbf{f})$

$$P_{\text{aug}}(s^0, \sigma^0 | s, \sigma, a) = \mathbb{1}_{f\sigma^0 = \alpha\sigma + \mathbf{f}(s, a, x)} g \sum_{i=1}^{M+1} z_i(\sigma) P_i(s^0 | s, a).$$

Finally, the augmented reward function r_{aug} satisfies

$$r_{\text{aug}}(s, \sigma, a) = \sum_{i=1}^{M+1} z_i(\sigma) r_i(s, a).$$

The augmented MDP \mathcal{M}_{aug} is closely related to the logistic DCMDP $(X, S, A, r, P, H, \mathbf{f}, \alpha)$. In fact, as we will show next, they both achieve the same optimal value. To see this, consider an MDP defined by the tuple $(S_1 \times S_2, A, P, r, H)$, and let $\phi : S_2 \rightarrow D$, where D is some known domain. Define the following set of deterministic policies

$$\Pi_{\text{aug}} = \{ \pi : S_1 \times S_2 \rightarrow A : \exists \eta : S_1 \times D \rightarrow [0, 1], \pi(s_1, s_2) = \eta(s_1, \phi(s_2)) g \}.$$

Define the augmented optimal value for some $s \in S_1 \times S_2$

$$V_{\text{aug},1}(s_1, s_2) = \max_{\pi \in \Pi_{\text{aug}}} \mathbb{E} \left[\sum_{t=1}^H r_t(s_t, a_t) \mid s_1 = s_1, s_2 = s_2, a_t = \pi_t(s_1, s_2) \right].$$

We apply the following proposition using the decomposition $S_1 = S$, and $S_2 = H$ as the set of possible trajectories in the known logistic DCMDP, where

$$\phi(\tau_h) := \sigma(\tau_h; \mathbf{f}) = \sum_{t=0}^{h-1} \alpha^{h-t-1} \mathbf{f}_t(s_t, a_t, x_t).$$

Proposition B.1 (Tennenholtz et al. (2022)). *Let $\mathcal{M} = (S_1 \times S_2, A, P, r, H)$. Assume for any $s_1, s_2 \in S_1 \times S_2, a \in A, P(s_1^0, \phi(s_2^0) | s_1, s_2, a) = P(s_1^0, \phi(s_2^0) | s_1, \phi(s_2), a)$ and $r(s_1, s_2, a) = g(s_1, a)$, for some deterministic function $g : S_1 \times A \rightarrow [0, 1]$. Then, for any $s_1, s_2 \in S_1 \times S_2$,*

$$V_{\text{aug},1}(s_1, s_2) = V_1(s_1, s_2).$$

This concludes our claim, proving that σ is indeed sufficient, as playing any policy in Π_{aug} achieves the same value.

B.2. Relation to TerMDPs

A special case of logistic DCMDPs are TerMDPs (Tennenholtz et al., 2022), which model exogenous, non-Markov termination in the environment. When terminated, the agent stops interacting with the environment and cannot collect additional rewards. This setup describes various real-world scenarios, such as passengers in autonomous vehicles or users abandoning a recommender systems. To model a TerMDP as a logistic DCMDP we let $X = \{0, 1\}$, $g = \mathbb{1}_{\text{term}}$, no term g , and define $r_i(s, a), P_i(s^0 | s, a)$ such that s_{term} is a sink state for which $r_i(s_{\text{term}}, a) = 0$. The reward in all other states is defined by $r_1(s, a)$. The transition probabilities are defined by $P_i(s^0 | s, a) = \begin{cases} s_{\text{term}}, & s = s_{\text{term}} - i = 0 \\ P_1(s^0 | s, a), & \text{o.w.} \end{cases}$.

TerMDPs use a cost functions $c_h(s_t, a_t)$ to define the probability of transitioning to the termination state, as $P(x_h = \text{term} | \tau_h) = z_0 \left(\sum_{t=1}^{h-1} c_t(s_t, a_t) \right)$ – a special case of logistic DCMDPs with a two-class, context-independent feature map, and a choice of $\alpha = 1$. Indeed, this choice of parameters defines a TerMDP as proposed in Tennenholtz et al. (2022).

Logistic DCMDPs let us consider generalized notions of such models, for which classes can reflect notions of trust, where humans become less susceptible to following recommendations from an agent, through situations where humans override an agent for short periods, to modeling the effects of changing moods.

C. Regret Analysis of LDC-UCB

In this section, we prove the regret bounds of Theorem 4.1. We start by defining the good event, which holds uniformly for all episode with probability $1 - \delta$. Then, we show that LDC-UCB is optimistic under the good event. Next, we decompose the regret to error terms of the reward, transition and latent features, and analyzing each of these terms result with the desired regret bounds.

We start by stating the bonuses which the algorithm uses:

$$b_{x,h}^{r,k}(s,a) = \min \left\{ \sqrt{\frac{\log \frac{8SAMHK}{\delta}}{n_h^k(s,a,x) - 1}}, 1 \right\}$$

$$b_{x,h}^{p,k}(s,a) = \min \left\{ H \sqrt{\frac{4S \log \frac{8SAMHK}{\delta}}{n_h^k(s,a,x) - 1}}, 2H \right\}$$

C.1. Failure Events

We define the following failure events.

$$F_k^r = \left\{ \exists s \in S, a \in A, x \in X, h \in [H] : |r_{x,h}(s,a) - \hat{r}_{x,h}^k(s,a)| > \min \left\{ \sqrt{\frac{\log \frac{2SAMHK}{\delta^\theta}}{n_h^k(s,a,x) - 1}}, 1 \right\} \right\}$$

$$F_k^p = \left\{ \exists s \in S, a \in A, x \in X, h \in [H] : \left\| P_{x,h}(j|s,a) - \hat{P}_{x,h}^k(j|s,a) \right\|_1 > \min \left\{ \sqrt{\frac{4S \log \frac{2SAMHK}{\delta^\theta}}{n_h^k(s,a,x) - 1}}, 2 \right\} \right\}$$

$$F_k^n = \left\{ \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[\frac{z_{i,h}^k}{\sqrt{n_h^k(s_h^k, a_h^k, i) - 1}} \mid F_k \right] > 18H^2 \log \left(\frac{1}{\delta^\theta} \right) + 2HS(M+1)A + 4\sqrt{H^2S(M+1)AK} \right\}$$

$$F_k^{\mathcal{F}, \text{global}} = \mathcal{F} \not\subseteq C_k(\delta^\theta),$$

where the definition of $C_k(\delta)$ can be found at Appendix G.

Then, we define the good event, where none of the aforementioned failure events ever occur, i.e.,

$$G = ([k \in [K]] \bar{F}_k^r) \setminus ([k \in [K]] \bar{F}_k^p) \setminus ([k \in [K]] \bar{F}_k^{\mathcal{F}, \text{global}}) \setminus \bar{F}_k^n$$

Lemma C.1. *Letting $\delta^\theta = \delta/4$, the event G holds with probability at least $1 - \delta$.*

Proof. We show that the probability that the events do not hold for all $k \in [K]$ is smaller than $\delta^\theta = \delta/4$.

Reward concentration. First observe that both the empirical and real rewards are bounded in $[0, 1]$, so if the minimizer in F_k^r is 1, the event never holds. Otherwise, for any fixed episode k , number of plays n , state s , action a , context x and timestep h , by Hoeffding's inequality, the estimation error is bounded w.p. $1 - \delta^\theta$ by $\sqrt{\frac{\log \frac{2}{\delta^\theta}}{n}}$. Taking the union bound over all possible values of $k, n - 1, s, a, x$ and h , w.p. at least $1 - \delta^\theta$, for all $k \in [K], s \in S, a \in A, x \in X, h \in [H]$, the estimation error is bounded by

$$|r_{x,h}(s,a) - \hat{r}_{x,h}^k(s,a)| \leq \sqrt{\frac{\log \frac{2SAMHK^2}{\delta^\theta}}{2n_h^k(s,a,x) - 1}} + \sqrt{\frac{\log \frac{2SAMHK}{\delta^\theta}}{n_h^k(s,a,x) - 1}}.$$

Finally, we remark that since $\delta^\theta = 1/4$, the event F_k^r never holds when $n_h^k(s,a,x) = 0$ since the bound is larger than 1.

In other words, $\Pr \{ [k \in [K]] \bar{F}_k^r \} \leq \delta^\theta$.

Transition concentration. By the exact same arguments as the reward concentration, while replacing Hoeffding's inequality by the concentration of the L_1 error of a probability estimator (Weissman et al., 2003), we also get $\Pr\{\sum_{k \geq [K]} \bar{F}_k^p\} \leq \delta^\ell$. Notice that the L_1 distance between any two probability distributions is bounded by 2, which justifies the minimization in the event.

Global feature estimation. By Lemma G.5, we have that

$$\Pr\left\{\sum_{k \geq [K]} \bar{F}_k^{\mathbf{f}, \text{global}}\right\} \leq \Pr\{f \notin C_k(\delta^\ell) \mid g \leq \delta^\ell\}.$$

Expected counts concentration. By Lemma E.7, we have that $\Pr\{\bar{F}^n\} \leq \delta^\ell$.

Fixing $\delta^\ell = \delta/4$ and taking the union bound concludes the proof. \square

C.2. Regret Analysis – Proof of Theorem 4.1

Theorem 4.1. Let $\lambda = \Theta(\frac{HM^{2.5}SA}{L})$. With probability at least $1 - \delta$, the regret of Algorithm 1 is

$$\text{Reg}(K) \leq \tilde{O}\left(\frac{P}{H^6 M^{4.5} S^2 A^2 L^2 \kappa K}\right).$$

Proof. Under the good event, the conditions of the regret decomposition lemma (Lemma E.6) hold with $\bar{r}, \hat{P}, Z^{\bar{F}_k}$ and $c = 4$ due to the truncated value iteration, truncated bonuses and value optimism lemma (Proposition C.5). Therefore, the regret can be decomposed in the following way

$$\begin{aligned} \text{Reg}(K) & \underbrace{\sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left| \bar{r}_{i,h}^k(s_h^k, a_h^k) - r_{i,h}(s_h^k, a_h^k) \right| \mid F_k \right]}_{(i)} \\ & + H \underbrace{\sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left\| \left(\hat{P}_h^k - P_h \right) (j_{s_h^k, a_h^k}^k) \right\|_1 \mid F_k \right]}_{(ii)} \\ & + 5H \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| Z_h^{\bar{F}_k} - Z_h^{\mathbf{f}} \right\|_1 \mid F_k \right]}_{(iii)}. \end{aligned}$$

By plugging in Lemma C.2, Lemma C.3 and Lemma C.4, which bound terms (i),(ii) and (iii), respectively, we get,

$$\begin{aligned} \text{Reg}(K) & \leq O\left(H^2 S \sqrt{MAK \log \frac{SAMHK}{\delta}}\right) \\ & + O\left(H^2 S \sqrt{MAK \log \frac{SAMHK}{\delta}}\right) \\ & + \tilde{O}\left(\frac{P}{S^2 A^2 H^6 M^{4.5} L^2 \kappa K}\right) \end{aligned}$$

Noticing that the last term is the dominant, we get

$$\text{Reg}(K) \leq \tilde{O}\left(\frac{P}{S^2 A^2 H^6 M^{4.5} L^2 \kappa K}\right),$$

which concludes the proof. \square

Now, we prove the lemmas that bounds the three terms in the regret decomposition in Theorem 4.1.

Lemma C.2 (Reward Concentration). *Under the good event, we have that:*

$$\begin{aligned}
 (i) &= \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left| \bar{r}_{i,h}^k(s_h^k, a_h^k) - r_{i,h}(s_h^k, a_h^k) \right| \mid F_{k-1} \right] \\
 &= 2H \sqrt{S \log \frac{8SAMHK}{\delta}} \left(18H^2 \log \left(\frac{4}{\delta} \right) + 2HS(M+1)A + 4\sqrt{H^2S(M+1)AK} \right) \quad (\text{Under } G) \\
 &= \mathcal{O} \left(H^2S \sqrt{MAK \log \frac{SAMHK}{\delta}} \right)
 \end{aligned}$$

Proof.

$$\begin{aligned}
 &\sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left| \bar{r}_{i,h}^k(s_h^k, a_h^k) - r_{i,h}(s_h^k, a_h^k) \right| \mid F_{k-1} \right] \\
 &\quad \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left(b_{i,h}^{r,k}(s_h^k, a_h^k) + b_{i,h}^{p,k}(s_h^k, a_h^k) + \left| \hat{r}_{i,h}^k(s_h^k, a_h^k) - r_{i,h}(s_h^k, a_h^k) \right| \right) \mid F_{k-1} \right] \\
 &\quad \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left(2b_{i,h}^{r,k}(s_h^k, a_h^k) + b_{i,h}^{p,k}(s_h^k, a_h^k) \right) \mid F_{k-1} \right] \quad (\text{Under } G) \\
 &\quad 2H \sqrt{S \log \frac{8SAMHK}{\delta}} \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[\frac{z_{i,h}^k}{\sqrt{n_h^k(s_h^k, a_h^k, i) - 1}} \mid F_{k-1} \right] \\
 &\quad 2H \sqrt{S \log \frac{8SAMHK}{\delta}} \left(18H^2 \log \left(\frac{4}{\delta} \right) + 2HS(M+1)A + 4\sqrt{H^2S(M+1)AK} \right) \quad (\text{Under } G) \\
 &= \mathcal{O} \left(H^2S \sqrt{MAK \log \frac{SAMHK}{\delta}} \right)
 \end{aligned}$$

□

Lemma C.3 (Transition Concentration). *Under the good event, we have that:*

$$\begin{aligned}
 (ii) &= H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left\| \left(\hat{P}_h^k - P_h \right) (j_{s_h^k, a_h^k}^k) \right\|_1 \mid F_{k-1} \right] \\
 &= \mathcal{O} \left(H^2S \sqrt{MAK \log \frac{SAMHK}{\delta}} \right)
 \end{aligned}$$

Proof.

$$\begin{aligned}
 &H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left\| \left(\hat{P}_h^k - P_h \right) (j_{s_h^k, a_h^k}^k) \right\|_1 \mid F_{k-1} \right] \\
 &\quad H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[z_{i,h}^k \frac{1}{H} b_{i,h}^{p,k}(s_h^k, a_h^k) \mid F_{k-1} \right] \quad (\text{Under } G) \\
 &\quad H \sqrt{4S \log \frac{8SAMHK}{\delta}} \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[\frac{z_{i,h}^k}{\sqrt{n_h^k(s_h^k, a_h^k, i) - 1}} \mid F_{k-1} \right] \\
 &\quad H \sqrt{4S \log \frac{8SAMHK}{\delta}} \left(18H^2 \log \left(\frac{4}{\delta} \right) + 2HS(M+1)A + 4\sqrt{H^2S(M+1)AK} \right) \quad (\text{Under } G) \\
 &= \mathcal{O} \left(H^2S \sqrt{MAK \log \frac{SAMHK}{\delta}} \right)
 \end{aligned}$$

□

Lemma C.4 (Latent Features Concentration). *Under the good event, if $L = \Omega(1)$ and $\lambda = \Theta(\frac{SAHM^{2.5}}{L})$, we have that:*

$$(iii) = 5H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| Z_h^{\bar{f}^k} - Z_h^{\mathbf{f}} \right\|_1 \mid F_{k-1} \right] \\ \tilde{O}\left(\frac{\rho}{S^2 A^2 H^6 M^{4.5} L^2 \kappa K}\right)$$

Proof.

$$\begin{aligned} & 5H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| Z_h^{\bar{f}^k} - Z_h^{\mathbf{f}} \right\|_1 \mid F_{k-1} \right] \\ & 5H \frac{\rho}{M+1} \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| Z_h^{\bar{f}^k} - Z_h^{\mathbf{f}} \right\|_2 \mid F_{k-1} \right] \\ & 10H \sqrt{(1+2L)(M+1)\kappa} \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\beta_k(4\delta) \left\| \mathbf{d}_h^k \right\|_{\mathbf{V}_{k-1}} \mid F_{k-1} \right] \quad (\text{Lemma G.8}) \\ & 10H \beta_K(4\delta) \sqrt{(1+2L)(M+1)\kappa} \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| \mathbf{d}_h^k \right\|_{\mathbf{V}_{k-1}} \mid F_{k-1} \right] \\ & 10H \beta_K(4\delta) \sqrt{(1+2L)(M+1)\kappa} \frac{\sqrt{2KH^3 MSA \log \frac{\kappa \lambda H SMA + k}{\kappa \lambda H SMA}}}{\max\{1, 1/\bar{\lambda}g\}} \quad (\text{Corollary E.9}) \\ & = O\left(\frac{\beta_K(4\delta)}{\max\{1, 1/\bar{\lambda}g\}} \sqrt{H^5 M^2 S A L \kappa K \log \frac{\kappa \lambda H SMA + k}{\kappa \lambda H SMA}}\right) \\ & = \tilde{O}\left(\frac{\rho}{S^2 A^2 H^6 M^{4.5} L^2 \kappa K}\right) \end{aligned}$$

For both the lemma and the corollary, we remind that $\mathbf{d}_h^k \leq 1$. For the last relation, recall that

$$\beta_k(\delta) = \frac{M^{3/2}(M+1)SAH}{\rho \bar{\lambda}} \left(\log \left(1 + \frac{k}{(M+1)SA\lambda} \right) + 2 \log \left(\frac{2}{\delta} \right) \right) + \sqrt{\frac{\lambda}{4M}} + \frac{\rho}{\bar{\lambda}L},$$

and assuming that $L = \Omega(1)$, we take $\lambda = \Theta(\frac{SAHM^{2.5}}{L})$, so $\beta_K(4\delta) = \tilde{O}(\frac{\rho}{SAHM^{2.5}L})$. □

C.3. Optimism in Logistic DCMDPs

In this section, we prove Proposition C.5, which allows us to apply the regret decomposition (Lemma E.6) necessary for proving Theorem 4.1.

We start by clearly stating the output value of the planning algorithm. For any $\mathbf{f} \in \mathcal{F}$, we define the truncated optimistic value under \mathbf{f} as the solution to the following value iteration problem:

$$\begin{aligned} \bar{V}_{H+1}^{k,\mathbf{f}}(s, \tau_{H+1}) &= 0, & \delta s &\in \mathcal{S}, \tau_{H+1} \in \mathcal{H}_{H+1} \\ \bar{V}_h^{k,\mathbf{f}}(s, \tau_h) &= \min \left\{ H, \max_a \left\{ \left[Z_h^{\mathbf{f}} \bar{r}_h^k \right](s, a, \tau_h) + \left[T_h^{\hat{P}_h^k, \mathbf{f}} \bar{V}_{h+1}^{k,\mathbf{f}} \right](s, a, \tau_h) \right\} \right\}, & \delta h &\in [H], s \in \mathcal{S}, \tau_h \in \mathcal{H}_h. \end{aligned}$$

Then, given an initial state s , we define the optimistic value of the DCMDP by $\bar{V}_1^k(s, \tau_h) = \max_{\mathbf{f} \in \mathcal{C}_k(\delta)} \bar{V}_1^{k,\mathbf{f}}(s, \tau_h)$. For this value, the following holds:

Proposition C.5. *Under the good event G , for any $k \in [K]$ and any initial state $s \in \mathcal{S}$, it holds that $\bar{V}_1^k(s) \leq V_1(s)$.*

Proof. Assume that G holds, and let $\bar{V}_h^{k,\mathbf{f}}(s, \tau_h)$, $\bar{V}_1^k(s, \tau_h)$ as defined by the beginning of the section. In the proof, we will show that for any $k \in [K]$, $s \in \mathcal{S}$, $h \in [H]$ and $\tau_h \in \mathcal{H}_h$ and $\mathbf{f} \in \mathcal{F}$, it holds that $\bar{V}_h^{k,\mathbf{f}}(s, \tau_h) \leq V_h^{k,\mathbf{f}}(s, \tau_h)$. Since under

G , we know that $\mathbf{f} \in \mathcal{C}_k(\delta)$, we then have that

$$\bar{V}_1^k(s) = \max_{\mathbf{f} \in \mathcal{C}_k(\delta)} \bar{V}_1^{k,\mathbf{f}}(s, \tau_h) \quad \max_{\mathbf{f} \in \mathcal{C}_k(\delta)} V_1^{\mathbf{f}}(s, \tau_h) \quad V_1(s),$$

which would conclude the prove. Throughout this proof, we assume w.l.o.g. that all optimistic values are smaller than H ; otherwise, they will be truncated to H , which still always optimistic since the rewards are in $[0, 1]$ and the horizon is H .

We prove that $\bar{V}_h^{k,\mathbf{f}}(s, \tau_h) = V_h^{\mathbf{f}}(s, \tau_h)$ by backward-induction. First notice that the claim holds for $h = H$, since

$$\begin{aligned} \bar{V}_H^{k,\mathbf{f}}(s, \tau_H; \mathbf{f}) = V_H^{\mathbf{f}}(s, \tau_H) &= \max_a \left\{ \left[Z_H^{\mathbf{f}} \bar{r}_H^k \right](s, a, \tau_H) \right\} = \max_a \left\{ \left[Z_H^{\mathbf{f}} r_H \right](s, a, \tau_H) \right\} \\ &\stackrel{(1)}{=} \left[Z_H^{\mathbf{f}} \bar{r}_H^k \right](s, a, \tau_H) = \left[Z_H^{\mathbf{f}} r_H \right](s, a, \tau_H) \\ &\quad \text{(for } a \in \arg \max_a \left\{ \left[Z_H^{\mathbf{f}} r_H \right](s, a, \tau_H) \right\}) \\ &= \left[Z_H^{\mathbf{f}} (\hat{r}_H^k + b_H^{r,k} - r) \right](s, a, \tau_H) \\ &= 0 \end{aligned} \tag{Lemma E.1}$$

Now let $h \in [H-1]$ and assume that the claim holds for $h+1$. Then, for

$$a \in \arg \max_a \left\{ \left[Z_h^{\mathbf{f}} r_h \right](s, a, \tau_h) + \left[T_h V_{h+1}^{\mathbf{f}} \right](s, a, \tau_h) \right\},$$

we have

$$\begin{aligned} \bar{V}_h^{k,\mathbf{f}}(s, \tau_h) = V_h^{\mathbf{f}}(s, \tau_h) &= \max_a \left\{ \left[Z_h^{\mathbf{f}} \bar{r}_h^k \right](s, a, \tau_h) + \left[T_h^{\hat{P}_h^k, \mathbf{f}} \bar{V}_{h+1}^{k,\mathbf{f}} \right](s, a, \tau_h) \right\} = \max_a \left\{ \left[Z_h^{\mathbf{f}} r_h \right](s, a, \tau_h) + \left[T_h V_{h+1}^{\mathbf{f}} \right](s, a, \tau_h) \right\} \\ &\stackrel{(1)}{=} \left[Z_h^{\mathbf{f}} \bar{r}_h^k \right](s, a, \tau_h) + \left[T_h^{\hat{P}_h^k, \mathbf{f}} \bar{V}_{h+1}^{k,\mathbf{f}} \right](s, a, \tau_h) = \left[Z_h^{\mathbf{f}} r_h \right](s, a, \tau_h) + \left[T_h V_{h+1}^{\mathbf{f}} \right](s, a, \tau_h) \\ &= \left[Z_h^{\mathbf{f}} \bar{r}_h^k \right](s, a, \tau_h) + \left[Z_h^{\mathbf{f}} r_h \right](s, a, \tau_h) + \left[\left(T_h^{\hat{P}_h^k, \mathbf{f}} - T_h \right) \bar{V}_{h+1}^{k,\mathbf{f}} \right](s, a, \tau_h) + \left[T_h \left(\bar{V}_{h+1}^{k,\mathbf{f}} - V_{h+1}^{\mathbf{f}} \right) \right](s, a, \tau_h) \\ &\stackrel{(2)}{=} \left[Z_h^{\mathbf{f}} \bar{r}_h^k \right](s, a, \tau_h) + \left[Z_h^{\mathbf{f}} r_h \right](s, a, \tau_h) + \left[\left(T_h^{\hat{P}_h^k, \mathbf{f}} - T_h \right) \bar{V}_{h+1}^{k,\mathbf{f}} \right](s, a, \tau_h) \\ &= \left[Z_h^{\mathbf{f}} (\bar{r}_h^k - r_h) \right](s, a, \tau_h) + \left[Z_h^{\mathbf{f}} \left(\hat{P}_h^k - P_h \right) \bar{V}_{h+1}^{k,\mathbf{f}} \right](s, a, \tau_h), \end{aligned}$$

where in (1) we used the definition of the max operator, and in (2) the induction step. Overall, replacing \bar{r}_h^k with its definition, we get that

$$\bar{V}_h^{k,\mathbf{f}}(s, \tau_h) - V_h^{\mathbf{f}}(s, \tau_h) = \left[Z_h^{\mathbf{f}} \left(\hat{r}_h^k + b_h^{r,k} - r_h \right) \right](s, a, \tau_h) + \left[Z_h^{\mathbf{f}} \left(\left(\hat{P}_h^k - P_h \right) \bar{V}_{h+1}^{k,\mathbf{f}} + b_h^{p,k} \right) \right](s, a, \tau_h) = 0,$$

where the second inequality is by Lemma E.1 and Lemma E.2, which hold under G .

□

D. Regret Analysis for Tractable LDC-UCB

In this section, we prove the regret bounds of Theorem 5.7. We start by defining the good event, which holds uniformly for all episode with probability $1 - \delta$. Then, we show that the Tractable LDC-UCB is optimistic under the good event. Next, we decompose the regret to error terms of the reward, transition and latent features, and analyzing each of these terms result with the desired regret bounds.

We start by stating the bonuses which the algorithm uses:

$$\begin{aligned} b_{x,h}^{r,k}(s,a) &= \min \left\{ \sqrt{\frac{\log \frac{8SAMHK}{\delta}}{n_h^k(s,a,x) - 1}}, 1 \right\} \\ b_{x,h}^{p,k}(s,a) &= \min \left\{ H \sqrt{\frac{4S \log \frac{8SAMHK}{\delta}}{n_h^k(s,a,x) - 1}}, 2H \right\} \\ b_{x,h}^{f,k}(s,a) &= \frac{2^{\rho} \bar{\kappa} \gamma_k(4\delta)}{\sqrt{n_h^k(s,a,x) + 4\lambda}} \end{aligned}$$

The confidence intervals can then be written as

$$\mathbf{C}_h^k = \left[\boldsymbol{\sigma}(\tau_h^k; \hat{\mathbf{f}}^k) \sum_{t=0}^{h-1} \alpha^{h-t} \mathbf{1}_{b_{x_t^k,t}^{f,k}(s_t^k, a_t^k)}, \boldsymbol{\sigma}(\tau_h^k; \hat{\mathbf{f}}^k) + \sum_{t=0}^{h-1} \alpha^{h-t} \mathbf{1}_{b_{x_t^k,t}^{f,k}(s_t^k, a_t^k)} \right]$$

D.1. Failure Events

We define the following failure events.

$$\begin{aligned} F_k^r &= \left\{ \exists s \in \mathcal{S}, a \in \mathcal{A}, x \in \mathcal{X}, h \in [H] : \hat{r}_{x,h}(s,a) - \hat{r}_{x,h}^k(s,a) > \min \left\{ \sqrt{\frac{\log \frac{2SAMHK}{\delta^\theta}}{n_h^k(s,a,x) - 1}}, 1 \right\} \right\} \\ F_k^p &= \left\{ \exists s \in \mathcal{S}, a \in \mathcal{A}, x \in \mathcal{X}, h \in [H] : \left\| P_{x,h}(j|s,a) - \hat{P}_{x,h}^k(j|s,a) \right\|_1 > \min \left\{ \sqrt{\frac{4S \log \frac{2SAMHK}{\delta^\theta}}{n_h^k(s,a,x) - 1}}, 2 \right\} \right\} \\ F_k^n &= \left\{ \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[\frac{z_{i,h}^k}{\sqrt{n_h^k(s_h^k, a_h^k, i) - 1}} \mid F_{k-1} \right] > 18H^2 \log \left(\frac{1}{\delta^\theta} \right) + 2HS(M+1)A + 4\sqrt{H^2 S(M+1)AK} \right\} \\ F_k^{f,\text{local}} &= \left\{ \exists s \in \mathcal{S}, a \in \mathcal{A}, x \in \mathcal{X}, i \in [M], h \in [H] : \left| \hat{f}_{i,h}^k(s,a,x) - f_{i,h}(s,a,x) \right| > \frac{2^{\rho} \bar{\kappa} \gamma_k(\delta^\theta)}{\sqrt{n_h^k(s,a,x) + 4\lambda}} \right\} \end{aligned}$$

where $\gamma_k(\delta)$ is defined in Proposition G.10.

Then, we define the good event, where none of the aforementioned failure events ever occur, i.e.,

$$G = ([_{k \in [K]} \bar{F}_k^r) \setminus ([_{k \in [K]} \bar{F}_k^p) \setminus ([_{k \in [K]} \bar{F}_k^{f,\text{local}}) \setminus \bar{F}^n$$

Lemma D.1. *Letting $\delta^\theta = \delta/4$, the event G holds with probability at least $1 - \delta$.*

Proof. The proof is almost identical to the one of Lemma C.1. We only need to prove that $\Pr \left\{ [_{k \in [K]} \bar{F}_k^{f,\text{local}} \right\} \leq \delta^\theta$, which directly follows by Lemma 5.2. \square

D.2. Regret Analysis – Proof of Theorem 5.7

Theorem 5.7. Let $\lambda = \Theta\left(\frac{HM^{2.5}SA}{L}\right)$. With probability at least $1 - \delta$, the regret of Algorithm 2 is

$$R(K) \quad \tilde{O}\left(\sqrt[{\rho}]{H^8 M^{6.5} S^2 A^2 L^4 \kappa K}\right).$$

Proof. Let

$$\bar{\sigma}_h^k(s, \tau_h) \triangleq \arg \max_{\bar{\sigma} \in \mathcal{C}_h^k} \max_a \left\{ \sum_{i=1}^{M+1} z_i(\bar{\sigma}) \bar{r}_{i,h}^k(s, a) + \sum_{i=1}^{M+1} z_i(\bar{\sigma}) \hat{P}_{i,h}^k(j_s, a)^T \bar{V}_{h+1}^k(\cdot, \tau_{h+1}) \right\},$$

and denote $\bar{Z}_h^k(s, \tau_h) = z(\bar{\sigma}_h^k(s, \tau_h))$.

Under the good event, the conditions of the regret decomposition lemma (Lemma E.6) hold with $\bar{r}, \hat{P}, \bar{Z}$ and $c = 4$ due to the truncated value iteration, truncated bonuses and value optimism lemma (Proposition C.5). Therefore, the regret can be decomposed in the following way,

$$\begin{aligned} \text{Reg}(K) & \underbrace{\sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left| r_{i,h}^k(s_h^k, a_h^k) - r_{i,h}(s_h^k, a_h^k) \right| \middle| F_{k-1} \right]}_{(i)} \\ & + H \underbrace{\sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left\| \left(\hat{P}_h^k - P_h \right) (j_{s_h^k}, a_h^k) \right\|_1 \middle| F_{k-1} \right]}_{(ii)} \\ & + 5H \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| \bar{Z}_h^k - Z_h^f \right\|_1 \middle| F_{k-1} \right]}_{(iii)}. \end{aligned}$$

The terms (i) and (iii) are identical to the ones in the proof of Theorem 4.1 (as the reward bonuses are identical), and thus can be bounded by Lemma C.2 and Lemma C.3. Term (ii) can be bounded by Lemma D.2.

Thus, we obtain,

$$\begin{aligned} \text{Reg}(K) & \quad O\left(H^2 S \sqrt{MAK \log \frac{SAMHK}{\delta}}\right) \\ & \quad + O\left(H^2 S \sqrt{MAK \log \frac{SAMHK}{\delta}}\right) \\ & \quad + \tilde{O}\left(\sqrt[{\rho}]{H^8 S^2 A^2 M^{6.5} L^4 \kappa K}\right) \\ & \quad \tilde{O}\left(\sqrt[{\rho}]{H^8 S^2 A^2 M^{6.5} L^4 \kappa K}\right) \end{aligned}$$

□

Lemma D.2 (Latent Features Concentration). *Under the good event, it holds that*

$$\begin{aligned} 5H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| Z_h^{\bar{f}^k} - Z_h^f \right\|_1 \middle| F_{k-1} \right] & \quad O\left(\left(\frac{1}{\lambda} - 1\right) \gamma_K(4\delta) \sqrt[{\rho}]{\kappa H^6 S M^3 AK}\right) \\ & \quad \tilde{O}\left(\sqrt[{\rho}]{H^8 S^2 A^2 M^{6.5} L^4 \kappa K}\right) \end{aligned}$$

Proof.

$$\begin{aligned}
 & 5H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| Z_h^{\bar{\mathbf{f}}^k} - Z_h^{\mathbf{f}} \right\|_1 \middle| F_{k-1} \right] \\
 &= 5H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| \mathbf{z}(\bar{\sigma}_h^k(s, \tau_h^k)) - \mathbf{z}(\sigma(\tau_h^k; \mathbf{f})) \right\|_1 \middle| F_{k-1} \right] \\
 &\stackrel{(1)}{=} 10H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[\left| z_i(\bar{\sigma}_h^k(s, \tau_h^k)) - z_i(\sigma(\tau_h^k; \mathbf{f})) \right| \middle| F_{k-1} \right] \\
 &\stackrel{(2)}{=} 2.5H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[\left| \bar{\sigma}_{i,h}^k(s, \tau_h^k) - \sigma_i(\tau_h^k; \mathbf{f}) \right| \middle| F_{k-1} \right]
 \end{aligned}$$

where relation (1) holds by substituting $z_{M+1}(\mathbf{x}) = 1 - \sum_{i=1}^M z_i(\mathbf{x})$ and applying the triangle inequality, and relation (2) is since the function $f(x) = e^x/(a + e^x)$ is $\frac{1}{4}$ -Lipschitz, and $z_i(\mathbf{x})$ can be represented as such a function of x_i . Then,

$$\begin{aligned}
 & 5H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| Z_h^{\bar{\mathbf{f}}^k} - Z_h^{\mathbf{f}} \right\|_1 \middle| F_{k-1} \right] \\
 & 2.5H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[\left| \bar{\sigma}_{i,h}^k(s, \tau_h^k) - \sigma_i(\tau_h^k; \mathbf{f}) \right| \middle| F_{k-1} \right] \\
 & 2.5H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[\left| \sigma_i(\tau_h; \hat{\mathbf{f}}^k) - \sigma_i(\tau_h^k; \mathbf{f}) \right| \middle| F_{k-1} \right] \\
 & + 2.5H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[\left| \sum_{t=0}^{h-1} \alpha^{h-t} {}^1b_{x_t^k, t}^{\mathbf{f}, k}(s_t^k, a_t^k) \right| \middle| F_{k-1} \right] \\
 & 2.5H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[\left| \sum_{t=0}^{h-1} \alpha^{h-t} {}^1b_{x_t^k, t}^{\mathbf{f}, k}(s_t^k, a_t^k) \right| \middle| F_{k-1} \right] \\
 & + 2.5H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^M \mathbb{E} \left[\left| \sum_{t=0}^{h-1} \alpha^{h-t} {}^1b_{x_t^k, t}^{\mathbf{f}, k}(s_t^k, a_t^k) \right| \middle| F_{k-1} \right] \tag{Under G} \\
 & = 5HM \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left| \sum_{t=0}^{h-1} \alpha^{h-t} {}^1b_{x_t^k, t}^{\mathbf{f}, k}(s_t^k, a_t^k) \right| \middle| F_{k-1} \right] \\
 & 5H^2M \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left| b_{x_h^k, h}^{\mathbf{f}, k}(s_h^k, a_h^k) \right| \middle| F_{k-1} \right] \\
 & 10H^2M \rho_{\kappa\gamma_K}^- \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k)} + 4\lambda} \middle| F_{k-1} \right] \\
 & 10H^2M \left(\frac{1}{2\lambda} - 1 \right) \rho_{\kappa\gamma_K}^- \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k)} - 1} \middle| F_{k-1} \right] \\
 & 10H^2M \left(\frac{1}{2\lambda} - 1 \right) \rho_{\kappa\gamma_K}^- \left(18H^2 \log \left(\frac{1}{\delta} \right) + 2HS(M+1)A + 4\sqrt{H^2S(M+1)AK} \right) \\
 & \tag{Lemma E.7} \\
 & = O \left(\left(\frac{1}{\lambda} - 1 \right) \gamma_K \rho_{\kappa H^6 S M^3 A K}^- \right)
 \end{aligned}$$

Now, note that by the definition in Proposition G.10, $\gamma_k(\delta) := (2 + 2L + \sqrt{2(1+L)})\beta_k(\delta) + \sqrt{\frac{2(1+L)HM}{\lambda}}\beta_k^2(\delta)$. Also, by Equation (12), $\beta_k(\delta) = \frac{M^{3/2}(M+1)SAH}{\rho_{\lambda}^-} \left(\log\left(1 + \frac{k}{(M+1)SA\lambda}\right) + 2\log\left(\frac{2}{\delta}\right) \right) + \sqrt{\frac{\lambda}{4M} + \rho_{\lambda}^-}L$. Plugging in these definitions we have that,

$$\begin{aligned} \gamma_k(\delta) &= (2 + 2L + \sqrt{2(1+L)}) \left(\frac{M^{3/2}(M+1)SAH}{\rho_{\lambda}^-} \left(\log\left(1 + \frac{k}{(M+1)SA\lambda}\right) + 2\log\left(\frac{2}{\delta}\right) \right) + \sqrt{\frac{\lambda}{4M} + \rho_{\lambda}^-}L \right) \\ &\quad + \sqrt{\frac{2(1+L)HM}{\lambda}} \left(\frac{M^{3/2}(M+1)SAH}{\rho_{\lambda}^-} \left(\log\left(1 + \frac{k}{(M+1)SA\lambda}\right) + 2\log\left(\frac{2}{\delta}\right) \right) + \sqrt{\frac{\lambda}{4M} + \rho_{\lambda}^-}L \right)^2 \\ &\quad + O\left(\frac{LM^{5/2}SAH}{\lambda} \log\left(\frac{1}{\delta} + \frac{k}{MSA\lambda\delta}\right) + L\sqrt{\frac{\lambda}{M} + \rho_{\lambda}^-}L^2 \right) \\ &\quad + \sqrt{\frac{2(1+L)HM}{\lambda}} O\left(\frac{M^5S^2AH^2}{\lambda} \log^2\left(\frac{1}{\delta} + \frac{k}{MSA\lambda\delta}\right) + \frac{\lambda}{M} + \lambda L^2 \right) \\ &= \tilde{O}\left(\lambda^{-1}LM^{5/2}SAH + \lambda^{1/2}LM^{-1/2} + \lambda^{1/2}L^2 \right) \\ &\quad + \tilde{O}\left(\lambda^{-3/2}L^{1/2}M^{11/2}S^2A^2H^{5/2} + \lambda^{1/2}L^{1/2}M^{-1/2}H^{1/2} + \lambda^{1/2}L^{5/2}M^{1/2}H^{1/2} \right) \\ &\quad + \tilde{O}\left(\lambda^{-1}LM^{5/2}SAH + \lambda^{1/2}LM^{-1/2} + \lambda^{-3/2}L^{1/2}M^{11/2}S^2A^2H^{5/2} \right. \\ &\quad \quad \left. + \lambda^{1/2}L^{1/2}M^{-1/2}H^{1/2} + \lambda^{1/2}L^{5/2}M^{1/2}H^{1/2} \right) \\ &\quad + \tilde{O}\left(L^2M^{7/4}S^{1/2}A^{1/2}H \right), \end{aligned}$$

where we used $\lambda = \frac{M^{2.5}SAH}{L}$ to minimize the above term.

Finally, plugging in this expression we have that

$$\begin{aligned} 5H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| Z_h^{\bar{f}^k} - Z_h^f \right\|_1 \mid F_{k-1} \right] &= O\left(\left(\frac{1}{\rho_{\lambda}^-} - 1 \right) \gamma_K(4\delta) \rho_{\lambda}^- \overline{\kappa H^6 S M^3 A K} \right) \\ &= \tilde{O}\left(L^2 M^{7/4} S^{1/2} A^{1/2} H \rho_{\lambda}^- \overline{\kappa H^6 S M^3 A K} \right) \\ &= \tilde{O}\left(\rho_{\lambda}^- \overline{H^8 S^2 A^2 M^{6.5} L^4 \kappa K} \right), \end{aligned}$$

where we assumed that $\lambda \geq 1$ □

D.3. Optimism in Tractable Logistic DCMDPs – Proof of Proposition 5.3

In this section, we prove Proposition 5.3, which allows us to apply the regret decomposition (Lemma E.6) necessary for proving Theorem 5.7.

Proposition 5.3 (Optimistic Value). *Let \bar{V}_h as defined in Equation (8). Then, w.h.p. $\bar{V}_1(s_1^k, \mathbf{C}_1^k) \leq V_1(s_1^k)$.*

Proof. We divide the proof into two steps. Defining $\bar{V}_h(s, \tau_h)$ the optimistic value function which follows the equations

$$\begin{aligned} \bar{V}_{H+1}^k(s, \tau_{H+1}) &= 0, \quad \forall s \in \mathcal{S}, \tau_{H+1} \in \mathcal{H}_{H+1}, \text{ and} \\ \bar{V}_h^k(s, \tau_h) &= \min \left\{ H, \max_a \max_{\bar{\sigma} \in \mathcal{C}(\tau_h)} \left\{ \sum_{i=1}^{M+1} z_i(\bar{\sigma}) \bar{r}_{i,h}^k(s, a) + \sum_{i=1}^{M+1} z_i(\bar{\sigma}) \hat{P}_{i,h}^k(j_s, a)^T \bar{V}_{h+1}^k(\cdot, \tau_{h+1}) \right\} \right\}, \\ \forall h \in [H], s \in \mathcal{S}, \tau_h \in \mathcal{H}_h, \end{aligned}$$

we first show that $\bar{V}_h^k(s, \mathbf{C}_h^k(\tau_h)) = \bar{V}_h^k(s, \tau_h)$ for all $k \in [K], h \in [H], s \in \mathcal{S}$ and $\tau_h \in \mathcal{H}_h$. This follows due to a simple induction; first notice that the claim trivially holds when $h = H + 1$, where both values are 0. Now fix $h \in [H]$ and assume

that $\bar{V}_{h+1}^k(s, \mathbf{C}_h^k(\tau_{h+1})) = \bar{V}_{h+1}^k(s, \tau_{h+1})$ for all $s \in S$ and $\tau_{h+1} \in H_{h+1}$. In the following, we prove that this implies $\bar{V}_h^k(s, \mathbf{C}_h^k(\tau_h)) = \bar{V}_h^k(s, \tau_h)$ for any $s \in S$ and $\tau_h \in H_h$, which prove the claim.

$$\begin{aligned}
 \bar{V}_h(s, \mathbf{C}_h^k(\tau_h)) &= \min \left\{ \max_{a \in \mathcal{A}} \max_{t \in \mathcal{T}(\bar{\mathbf{Q}})} \sum_{i=0}^M z_i \left(\mathbf{th}_t \left(\bar{\mathbf{Q}}, \mathbf{C}_h^k(\tau_h) \right) \right) \bar{Q}_i(s, a, \mathbf{C}_h^k(\tau_h)), H \right\} \\
 &\stackrel{(1)}{=} \min \left\{ \max_{a \in \mathcal{A}} \max_{\bar{\boldsymbol{\sigma}} \in \mathcal{C}_h^k(\tau_h)} \sum_{i=0}^M z_i(\bar{\boldsymbol{\sigma}}) \bar{Q}_i(s, a, \mathbf{C}_h^k(\tau_h)), H \right\} \\
 &= \min \left\{ \max_{a \in \mathcal{A}} \max_{\bar{\boldsymbol{\sigma}} \in \mathcal{C}_h^k(\tau_h)} \sum_{i=0}^M z_i(\bar{\boldsymbol{\sigma}}) \left(\bar{r}_i(s, a) + \mathbb{E}_{s^\theta} \hat{P}_i(j_s, a) [\bar{V}_{h+1}(s^\theta, \mathbf{C}_{h+1}(a, i))] \right), H \right\} \\
 &\stackrel{(2)}{=} \min \left\{ \max_{a \in \mathcal{A}} \max_{\bar{\boldsymbol{\sigma}} \in \mathcal{C}_h^k(\tau_h)} \sum_{i=0}^M z_i(\bar{\boldsymbol{\sigma}}) \left(\bar{r}_i(s, a) + \mathbb{E}_{s^\theta} \hat{P}_i(j_s, a) [\bar{V}_{h+1}(s^\theta, \mathbf{C}_{h+1}^k)] \right), H \right\} \\
 &\stackrel{(3)}{=} \min \left\{ \max_{a \in \mathcal{A}} \max_{\bar{\boldsymbol{\sigma}} \in \mathcal{C}_h^k(\tau_h)} \sum_{i=0}^M z_i(\bar{\boldsymbol{\sigma}}) \left(\bar{r}_i(s, a) + \mathbb{E}_{s^\theta} \hat{P}_i(j_s, a) [\bar{V}_{h+1}^k(s, \tau_h)] \right), H \right\} \\
 &= \bar{V}_h^k(s, \tau_h)
 \end{aligned}$$

Relation (1) is by Lemma 5.6, which proves that when the confidence interval of a multinomial function is rectangular, one of the maximizers of an linear combination w.r.t. this function is a threshold function; therefore, the maximum over threshold functions achieves the same value at the rectangular set. Relation (2) is by the definition of $\mathbf{C}_{h+1}(a_h, x_h)$ at Algorithm 4, which implies that $\mathbf{C}_{h+1}(a_h, x_h) = \mathbf{C}_{h+1}^k$. Finally, (3) is by the induction hypothesis.

Next, we prove that under the good event, $\bar{V}_h^k(s, \tau_h) = V_h(s, \tau_h)$ for all $k \in [K], h \in [H], s \in S$ and $\tau_h \in H_h$. This claim is also proved by induction and clearly holds when $h = H + 1$, when all values equal zero. Assume that $\bar{V}_{h+1}^k(s, \tau_{h+1}) = V_{h+1}(s, \tau_{h+1})$ for all $s \in S, \tau_{h+1} \in H_{h+1}$. Also, assume w.l.o.g. that $\bar{V}_h^k(s, \tau_h) < H$, otherwise the claim trivially holds. Then, denoting

$$a \in \arg \max_a \left\{ \left[Z_h^{\mathbf{f}} r_h \right] (s, a, \tau_h) + \left[T_h V_{h+1}^{\mathbf{f}} \right] (s, a, \tau_h) \right\},$$

and under G ,

$$\begin{aligned}
 &\bar{V}_h^k(s, \tau_h) - V_h(s, \tau_h) \\
 &= \max_a \max_{\bar{\boldsymbol{\sigma}} \in \mathcal{C}_h^k} \left\{ \sum_{i=1}^{M+1} z_i(\bar{\boldsymbol{\sigma}}) \bar{r}_{i,h}^k(s, a) + \sum_{i=1}^{M+1} z_i(\bar{\boldsymbol{\sigma}}) \hat{P}_{i,h}^k(j_s, a)^T \bar{V}_{h+1}^k(s, \tau_{h+1}) \right\} \\
 &\quad \max_a \left\{ \left[Z_h^{\mathbf{f}} r_h \right] (s, a, \tau_h) + \left[T_h V_{h+1}^{\mathbf{f}} \right] (s, a, \tau_h) \right\} \\
 &\stackrel{(1)}{=} \max_{\bar{\boldsymbol{\sigma}} \in \mathcal{C}_h^k} \left\{ \sum_{i=1}^{M+1} z_i(\bar{\boldsymbol{\sigma}}) \bar{r}_{i,h}^k(s, a) + \sum_{i=1}^{M+1} z_i(\bar{\boldsymbol{\sigma}}) \hat{P}_{i,h}^k(j_s, a)^T \bar{V}_{h+1}^k(s, \tau_{h+1}) \right\} \\
 &\quad \left[Z_h^{\mathbf{f}} r_h \right] (s, a, \tau_h) - \left[T_h V_{h+1}^{\mathbf{f}} \right] (s, a, \tau_h) \\
 &\stackrel{(2)}{=} \left[Z_h^{\mathbf{f}} \bar{r}_h^k \right] (s, a, \tau_h) - \left[T_h^{\hat{P}_h^k, \mathbf{f}} \bar{V}_{h+1}^k \right] (s, a, \tau_h) - \left[Z_h^{\mathbf{f}} r_h \right] (s, a, \tau_h) + \left[T_h V_{h+1}^{\mathbf{f}} \right] (s, a, \tau_h) \\
 &= \left[Z_h^{\mathbf{f}} \bar{r}_h^k \right] (s, a, \tau_h) - \left[Z_h^{\mathbf{f}} r_h \right] (s, a, \tau_h) + \left[\left(T_h^{\hat{P}_h^k, \mathbf{f}} - T_h \right) \bar{V}_{h+1}^k \right] (s, a, \tau_h) + \left[T_h \left(\bar{V}_{h+1}^k - V_{h+1}^{\mathbf{f}} \right) \right] (s, a, \tau_h) \\
 &\stackrel{(3)}{=} \left[Z_h^{\mathbf{f}} \bar{r}_h^k \right] (s, a, \tau_h) - \left[Z_h^{\mathbf{f}} r_h \right] (s, a, \tau_h) + \left[\left(T_h^{\hat{P}_h^k, \mathbf{f}} - T_h \right) \bar{V}_{h+1}^k \right] (s, a, \tau_h) \\
 &= \left[Z_h^{\mathbf{f}} \left(\bar{r}_h^k - r_h \right) \right] (s, a, \tau_h) + \left[Z_h^{\mathbf{f}} \left(\hat{P}_h^k - P_h \right) \bar{V}_{h+1}^k \right] (s, a, \tau_h).
 \end{aligned}$$

In (1) we used the definition of the max operator. Relation (2) holds since under the good event,

$$\begin{aligned}
 C(\sigma(\tau_h^k; \mathbf{f})) &= \sum_{t=0}^{h-1} \alpha^{h-t-1} \mathbf{f}_t(s_t, a_t, x_t) \\
 &\geq \left[\sigma(\tau_h; \hat{\mathbf{f}}^k) \sum_{t=0}^{h-1} \alpha^{h-t-1} b_{x_t^k, t}^{\mathbf{f}, k}(s_t^k, a_t^k), \sigma(\tau; \hat{\mathbf{f}}^k) + \sum_{t=0}^{h-1} \alpha^{h-t-1} b_{x_t^k, t}^{\mathbf{f}, k}(s_t^k, a_t^k) \right] \quad (\text{Under } G) \\
 &= C_h^k.
 \end{aligned}$$

In (3), we used the induction step the induction step. Overall, replacing \bar{r}_h^k with its definition, we get that

$$\begin{aligned}
 \bar{V}_h^k(s, \tau_h) - V_h(s, \tau_h) &\leq \left[Z_h^{\mathbf{f}} \left(\hat{r}_h^k + b_h^{r, k} - r_h \right) \right](s, a, \tau_h) + \left[Z_h^{\mathbf{f}} \left(\left(\hat{P}_h^k - P_h \right) \bar{V}_{h+1}^{\mathbf{f}, k} + b_h^{p, k} \right) \right](s, a, \tau_h) \\
 &\leq 0,
 \end{aligned}$$

where the second inequality is by Lemma E.1 and Lemma E.2, which hold under G . □

E. Useful Lemmas

E.1. Optimism Lemmas

Lemma E.1 (Reward Optimism). *For any $k \geq 1$, define the event*

$$F_k^r = \left\{ \exists s \in S, a \in A, i \in [M], h \in [H] : j r_{i,h}(s, a) - \hat{r}_{i,h}^{k,r}(s, a) > b_{i,h}^k(s, a) \right\}.$$

Then, under \bar{F}_k^r , for any $f \in F$, $h \in [H]$, $s \in S$, $a \in A$ and $\tau_h \in H_h$, it holds that

$$\left[Z_h^f(\hat{r}_h^k + b_h^{r,k} - r) \right](s, a, \tau_h) \leq 0$$

Proof. The result directly follows by the definition of \bar{F}_k^r , since

$$\begin{aligned} \left[Z_h^f(\hat{r}_h^k + b_h^{r,k} - r) \right](s, a, \tau_h) &= \min_i \left\{ (\hat{r}_{i,h}^k(s, a) - r_i(s, a)) + b_{i,h}^{r,k}(s, a) \right\} \\ &= \min_i \left\{ b_{i,h}^{r,k}(s, a) + b_{i,h}^{r,k}(s, a) \right\} \quad (\text{Under } \bar{F}_k^r) \\ &= 0 \end{aligned}$$

□

Lemma E.2 (Transition Optimism). *For any $k \geq 1$, define the event*

$$F_k^p = \left\{ \exists s \in S, a \in A, i \in [M], h \in [H] : \left\| P_{i,h}(\cdot | j, s, a) - \hat{P}_{i,h}^k(\cdot | j, s, a) \right\|_1 > \frac{1}{H} b_{i,h}^{p,k}(s, a) \right\}.$$

Then, under \bar{F}_k^p , for any $f \in F$, $h \in [H]$, $s \in S$, $a \in A$, $\tau_h \in H_h$ and $V \in [0, H]^S$, it holds that

$$\left[Z_h^f \left(\left(\hat{P}_h^k - P_h \right) V + b_h^{p,k} \right) \right](s, a, \tau_h) \leq 0$$

Proof. The result directly follows by the definition of \bar{F}_k^p and Cauchy-Schwartz inequality, since

$$\begin{aligned} \left[Z_h^f \left(\left(\hat{P}_h^k - P_h \right) V + b_h^{p,k} \right) \right](s, a, \tau_h) &= \min_i \left\{ \left[\left(\hat{P}_h^k - P_h \right) V \right](s, a) + b_{i,h}^{p,k}(s, a) \right\} \\ &= \min_i \left\{ \left\| P_{i,h}(\cdot | j, s, a) - \hat{P}_{i,h}^k(\cdot | j, s, a) \right\|_1 k V k_\gamma + b_{i,h}^{p,k}(s, a) \right\} \quad (\text{C.S}) \\ &= \min_i \left\{ \frac{1}{H} b_{i,h}^{p,k}(s, a) H + b_{i,h}^{p,k}(s, a) \right\} \quad (\text{Under } \bar{F}_k^p) \\ &= 0 \end{aligned}$$

□

E.2. Decomposition Lemmas

Lemma E.3.

$$\begin{aligned} V &: S \rightarrow H \times \mathbb{R}, \\ Z^{(1)}, Z^{(2)} &: S \times A \rightarrow H \times \Delta_X, \\ r^{(1)}, r^{(2)} &: S \times A \rightarrow H \times \mathbb{R}, \text{ and} \\ P^{(1)}, P^{(2)} &: S \times A \rightarrow H \times \Delta_S. \end{aligned}$$

Then, for any $s \in S, a \in A, h \in [H], \tau_h \in H$

$$\begin{aligned} & \left[Z_h^{(1)} r_h^{(1)} + T_h^{P^{(1)}, Z^{(1)}} V_{h+1} \right] (s, a, \tau_h) - \left[Z_h^{(2)} r_h^{(2)} + T_h^{P^{(2)}, Z^{(2)}} V_{h+1} \right] (s, a, \tau_h) \\ &= \left[Z_h^{(2)} \begin{pmatrix} r_h^{(1)} & r_h^{(2)} \end{pmatrix} \right] (s, a, \tau_h) \\ &+ \left[\begin{pmatrix} Z_h^{(1)} & Z_h^{(2)} \end{pmatrix} \begin{pmatrix} r_h^{(1)} + P_h^{(1)} V_{h+1} \\ r_h^{(2)} + P_h^{(2)} V_{h+1} \end{pmatrix} \right] (s, a, \tau_h) \\ &+ \left[Z_h^{(2)} \begin{pmatrix} P_h^{(1)} & P_h^{(2)} \end{pmatrix} V_{h+1} \right] (s, a, \tau_h). \end{aligned}$$

Proof. We have that

$$\begin{aligned} & \left[Z_h^{(1)} r_h^{(1)} + T_h^{P^{(1)}, Z^{(1)}} V_{h+1} \right] (s, a, \tau_h) - \left[Z_h^{(2)} r_h^{(2)} + T_h^{P^{(2)}, Z^{(2)}} V_{h+1} \right] (s, a, \tau_h) \\ &= \left[\begin{pmatrix} Z_h^{(1)} & Z_h^{(2)} \end{pmatrix} r_h^{(1)} \right] (s, a, \tau_h) + \left[Z^{(2)} \begin{pmatrix} r_h^{(1)} & r_h^{(2)} \end{pmatrix} \right] (s, a, \tau_h) + \left[\begin{pmatrix} T_h^{P^{(1)}, Z^{(1)}} & T_h^{P^{(2)}, Z^{(2)}} \end{pmatrix} V_{h+1} \right] (s, a, \tau_h) \\ &= \left[\begin{pmatrix} Z_h^{(1)} & Z_h^{(2)} \end{pmatrix} r_h^{(1)} \right] (s, a, \tau_h) + \left[Z^{(2)} \begin{pmatrix} r_h^{(1)} & r_h^{(2)} \end{pmatrix} \right] (s, a, \tau_h) + \left[\begin{pmatrix} Z_h^{(1)} P_h^{(1)} & Z_h^{(2)} P_h^{(2)} \end{pmatrix} V_{h+1} \right] (s, a, \tau_h) \\ &= \left[\begin{pmatrix} Z_h^{(1)} & Z_h^{(2)} \end{pmatrix} r_h^{(1)} \right] (s, a, \tau_h) + \left[Z^{(2)} \begin{pmatrix} r_h^{(1)} & r_h^{(2)} \end{pmatrix} \right] (s, a, \tau_h) \\ &+ \left[\begin{pmatrix} Z_h^{(1)} & Z_h^{(2)} \end{pmatrix} P_h^{(1)} V_{h+1} \right] (s, a, \tau_h) + \left[Z_h^{(2)} \begin{pmatrix} P_h^{(1)} & P_h^{(2)} \end{pmatrix} V_{h+1} \right] (s, a, \tau_h) \\ &= \left[Z_h^{(2)} \begin{pmatrix} r_h^{(1)} & r_h^{(2)} \end{pmatrix} \right] (s, a, \tau_h) + \left[\begin{pmatrix} Z_h^{(1)} & Z_h^{(2)} \end{pmatrix} \begin{pmatrix} r_h^{(1)} + P_h^{(1)} V_{h+1} \\ r_h^{(2)} + P_h^{(2)} V_{h+1} \end{pmatrix} \right] (s, a, \tau_h) + \left[Z_h^{(2)} \begin{pmatrix} P_h^{(1)} & P_h^{(2)} \end{pmatrix} V_{h+1} \right] (s, a, \tau_h). \end{aligned}$$

This completes the proof. \square

Next, recall that by embedding the history into the state, every DCMDP can be represented as an MDP. This equivalence will allow us to apply the following lemma on DCMDPs:

Lemma E.4 (Value difference lemma, e.g., Dann et al. (2017), Lemma E.15). *Consider two MDPs $\mathcal{M} = (S, A, P, r, H)$ and $\mathcal{M}^\theta = (S, A, P^\theta, r^\theta, H)$. For any policy π and any s, h , the following relation holds:*

$$\begin{aligned} & V_h^\pi(s; \mathcal{M}^\theta) - V_h^\pi(s; \mathcal{M}) \\ &= \mathbb{E} \left[\sum_{t=h}^H \begin{pmatrix} r_t^\theta(s_t, a_t) & r_t(s_t, a_t) \end{pmatrix} + \begin{pmatrix} P^\theta & P \end{pmatrix} \begin{pmatrix} j_{s_t, a_t}^T V_{t+1}^\pi(\cdot; \mathcal{M}^\theta) \\ j_{s_t, a_t}^T V_{t+1}^\pi(\cdot; \mathcal{M}) \end{pmatrix} \middle| s_h = s, \pi, P \right] \end{aligned}$$

Corollary E.5 (Truncated value difference lemma). *Consider two MDPs $\mathcal{M} = (S, A, P, r, H)$ and $\mathcal{M}^\theta = (S, A, P^\theta, r^\theta, H)$. Also, for any $C \in \mathbb{R}$, define the truncated value of a policy π under MDP \mathcal{M} by the solution to the truncated dynamic programming problem*

$$\begin{aligned} & V_{H+1}^\pi(s; \mathcal{M}, C) = 0, & \forall s \in S \\ & V_h^\pi(s; \mathcal{M}, C) = \mathbb{E}_a \left[\min \{ C, r_h(s, a) + P(j_{s, a})^T V_{h+1}^\pi(\cdot; \mathcal{M}, C) \} \right], & \forall h \in [H], s \in S. \end{aligned}$$

Then, for any policy π , any $s \in S, h \in [H]$ and any $C \in \mathbb{R}$, the following relation holds:

$$\begin{aligned} & V_h^\pi(s; \mathcal{M}^\theta, C) - V_h^\pi(s; \mathcal{M}) \\ &= \mathbb{E} \left[\sum_{t=h}^H \begin{pmatrix} r_t^\theta(s_t, a_t) & r_t(s_t, a_t) \end{pmatrix} + \begin{pmatrix} P^\theta & P \end{pmatrix} \begin{pmatrix} j_{s_t, a_t}^T V_{t+1}^\pi(\cdot; \mathcal{M}^\theta, C) \\ j_{s_t, a_t}^T V_{t+1}^\pi(\cdot; \mathcal{M}, C) \end{pmatrix} \middle| s_h = s, \pi, P \right] \end{aligned}$$

Proof. We build an MDP whose value (without truncation) is $V_h^\pi(s; \mathcal{M}^\theta, C)$ and its reward are always smaller than the rewards of \mathcal{M}^θ . In particular, for any $h \in [H], s \in S$ and $a \in A$, define the new reward function

$$\bar{r}_h(s, a) = r_h^\theta(s, a) - \max \{ 0, r_h^\theta(s, a) + P^\theta(j_{s, a})^T V_{h+1}^\pi(s; \mathcal{M}^\theta, C) - C \} - r_h^\theta(s, a), \quad (10)$$

and denote $\bar{\mathcal{M}} = (S, A, P^\theta, \bar{r}, H)$. Clearly, $V_{H+1}^\pi(s; \bar{\mathcal{M}}) = V_{H+1}^\pi(s; \mathcal{M}^\theta, C) = 0$. Now assume by induction the equality holds for all $t > h$ and all $s \in S$.

Let $s \in S$ be some state. If for some $a \in A$, the maximizer in Equation (10) equals zero, then there was no truncation in the value iteration and so, by the induction hypothesis, we get

$$\begin{aligned} \bar{r}_h(s, a) + P^\theta(j, s, a)^T V_{h+1}^\pi(\cdot; \bar{\mathcal{M}}) &= r_h^\theta(s, a) + P^\theta(j, s, a)^T V_{h+1}^\pi(\cdot; \mathcal{M}^\theta, C) \\ &= \min\{C, r_h(s, a) + P(j, s, a)^T V_{h+1}^\pi(\cdot; \mathcal{M}^\theta, C)\}. \end{aligned}$$

On the other hand, if the maximizer in Equation (10) is not zero, then one can easily verify that

$$\bar{r}_h(s, a) + P^\theta(j, s, a)^T V_{h+1}^\pi(\cdot; \bar{\mathcal{M}}) = \min\{C, r_h(s, a) + P^\theta(j, s, a)^T V_{h+1}^\pi(\cdot; \mathcal{M}^\theta, C)\} = C$$

Therefore, this equality holds for all $a \in A$ and thus

$$\begin{aligned} V_h^\pi(s; \bar{\mathcal{M}}) &= \mathbb{E}_a \pi [\bar{r}_h(s, a) + P^\theta(j, s, a)^T V_{h+1}^\pi(\cdot; \bar{\mathcal{M}})] \\ &= \mathbb{E}_a \pi [\min\{C, r_h(s, a) + P^\theta(j, s, a)^T V_{h+1}^\pi(\cdot; \mathcal{M}^\theta, C)\}] \\ &= V_h^\pi(s; \mathcal{M}^\theta, C), \end{aligned}$$

and by induction, this equality holds for all $h \in [H]$ and $s \in S$. Now, using this fact with Lemma E.4 on \mathcal{M} and $\bar{\mathcal{M}}$, we get:

$$\begin{aligned} &V_h^\pi(s; \mathcal{M}^\theta, C) - V_h^\pi(s; \mathcal{M}) \\ &= \mathbb{E} \left[\sum_{t=h}^H (\bar{r}_t(s_t, a_t) - r_t(s_t, a_t)) + (P^\theta - P)(j, s_t, a_t)^T V_{t+1}^\pi(\cdot; \mathcal{M}^\theta, C) \middle| s_h = s, \pi, P \right] \\ &= \mathbb{E} \left[\sum_{t=h}^H (r_t^\theta(s_t, a_t) - r_t(s_t, a_t)) + (P^\theta - P)(j, s_t, a_t)^T V_{t+1}^\pi(\cdot; \mathcal{M}^\theta, C) \middle| s_h = s, \pi, P \right] \end{aligned}$$

where the inequality is since $\bar{r}_h(s, a) \geq r_h^\theta(s, a)$ for all h, s, a . \square

We are now ready to present the general regret decomposition lemma.

Lemma E.6 (Regret Decomposition). *Assume that there exist an optimistic value function \bar{V}_h^k such that the following hold:*

1. **Value representation.** For all $k \in [K]$ and $h \in [H]$, there exist $\bar{Z}_h^k : S \times A \times H_h \times \mathcal{V} \times \Delta_X$, $\bar{r}_h^k : S \times A \times \mathcal{V} \times \mathbb{R}$ and $\bar{P}_h^k : S \times A \times \mathcal{V} \times \Delta_S$ such that for all $k = 1, h \in [H]$, $s \in S$ and $\tau_h \in H_h$, it holds that

$$\bar{V}_h^k(s, \tau_h) = \bar{Z}_h^k \bar{r}_h^k + T_h^{\bar{P}_h^k, \bar{J}_h^k} (j, s_h^k, a_h^k, \tau_h^k)^T \bar{V}_{h+1}^k(\cdot, \tau_{h+1}^k).$$

2. **Boundedness.** For all $k = 1, h \in [H]$, $s \in S$, $a \in A$ and $i \in [M+1]$ and $\tau_h \in H_h$, it holds that $0 \leq \bar{V}_h^k(s, \tau_h) \leq H$ and $0 \leq \bar{r}_{i,h}^k(s, a, \tau_h) \leq cH$ for some $c > 0$.
3. **Optimism.** For all $k = 1$, it holds that $\bar{V}_h^k(s_1^k) \geq V_1(s_1^k)$.

Then, the regret can be bounded by

$$\begin{aligned} \text{Reg}(K) &= \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} [z_{i,h}^k | \bar{r}_{i,h}^k(s_h^k, a_h^k) - r_{i,h}(s_h^k, a_h^k)| \mid F_{k-1}] \\ &+ H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} [z_{i,h}^k \| (\bar{P}_h^k - P_h)(j, s_h^k, a_h^k) \|_1 \mid F_{k-1}] \\ &+ (c+1)H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} [\| \bar{Z}_h^k - Z_h^f \|_1 \mid F_{k-1}] \end{aligned}$$

Proof.

$$\begin{aligned}
 \text{Reg}(K) &= \sum_{k=1}^K V_1(s_1^k) - V_1^{\pi^k}(s_1^k) \\
 &\quad \sum_{k=1}^K \bar{V}_1^k(s_1^k) - V_1^{\pi^k}(s_1^k) \tag{Optimism} \\
 &\quad \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left(\bar{Z}_h^k \bar{r}_h^k \quad Z_h^f r_h \right) (s_h^k, a_h^k, \tau_h^k) + (T_h^{\bar{P}_h^k, \bar{J}_h^k} - T_h) (js_h^k, a_h^k, \tau_h^k)^T \bar{V}_{h+1}^k(\cdot, \tau_{h+1}^k) \mid F_{k-1} \right] \\
 &\tag{Corollary E.5} \\
 &\quad \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left[Z_h^f \left(\bar{r}_h^k \quad r_h \right) \right] (s_h^k, a_h^k, \tau_h^k) \mid F_{k-1} \right]}_{(i)} \\
 &\quad + \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left(\bar{Z}_h^k \quad Z_h^f \right) \left(\bar{r}_h^k(s_h^k, a_h^k) + \bar{P}_h^k(js_h^k, a_h^k) \bar{V}_{h+1}^k(\cdot, \tau_{h+1}^k) \right) \mid F_{k-1} \right]}_{(ii)} \\
 &\quad + \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[Z_h^f \left[\bar{P}_h^k - P_h \right] (js_h^k, a_h^k, \tau_h^k)^T \bar{V}_{h+1}^k(\cdot, \tau_{h+1}^k) \mid F_{k-1} \right]}_{(iii)} \tag{Lemma E.3}
 \end{aligned}$$

Notice that in the application of Lemma E.4, which was applied w.r.t. π^k , we used the fact that any DCMDP can be represented as an MDP whose history was embedded into the state. We now bound each of the terms of the decomposition.

Reward error

$$\begin{aligned}
 (i) &= \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\sum_{i=1}^{M+1} z_{i,h}^k \left(\bar{r}_{i,h}^k(s_h^k, a_h^k) - r_{i,h}(s_h^k, a_h^k) \right) \mid F_{k-1} \right] \\
 &\quad \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left| \bar{r}_{i,h}^k(s_h^k, a_h^k) - r_{i,h}(s_h^k, a_h^k) \right| \mid F_{k-1} \right]
 \end{aligned}$$

Latent features error

$$\begin{aligned}
 (ii) &\quad \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| \bar{Z}_h^k \quad Z_h^f \right\|_1 \left\| \bar{r}_h^k(s_h^k, a_h^k) + \hat{P}_h^k(js_h^k, a_h^k)^T \bar{V}_{h+1}^k(\cdot, \tau_{h+1}^k) \right\|_\gamma \mid F_{k-1} \right] \tag{Hölder} \\
 &\quad (c+1)H \frac{\rho}{M+1} \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| \bar{Z}_h^k \quad Z_h^f \right\|_1 \mid F_{k-1} \right]
 \end{aligned}$$

where the last inequality is since the optimistic value is bounded in $[0, H]$ and the reward is in $[0, cH]$.

Transition error

$$\begin{aligned}
 (iii) &= \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\sum_{i=1}^{M+1} z_{i,h}^k \left(\left(\hat{P}_h^k - P_h \right) (js_h^k, a_h^k)^T \bar{V}_{h+1}^k(\cdot, \tau_{h+1}^k) \right) \mid F_{k-1} \right] \\
 &\quad \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left\| \left(\hat{P}_h^k - P_h \right) (js_h^k, a_h^k) \right\|_1 \left\| \bar{V}_{h+1}^k(\cdot, \tau_{h+1}^k) \right\|_\gamma \mid F_{k-1} \right] \tag{Hölder} \\
 &\quad H \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[z_{i,h}^k \left\| \left(\hat{P}_h^k - P_h \right) (js_h^k, a_h^k) \right\|_1 \mid F_{k-1} \right] \tag{Boundedness}
 \end{aligned}$$

Combining all bounds concludes the proof. \square

E.3. Visitation-Summation Lemmas

Lemma E.7 (Expected Cumulative Visitation Bound, Lemma 22, Efroni et al. (2020), adapted to DCMDPs). *Let $\mathcal{F}_k, \mathcal{G}_{k=1}^K$ be the natural filtration. Then, with probability greater than $1 - \delta$ it holds that*

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[\frac{z_{i,h}^k}{\sqrt{n_h^k(s_h^k, a_h^k, i) - 1}} \middle| \mathcal{F}_{k-1} \right] &= \sum_{k=1}^K \mathbb{E} \left[\sum_{h=1}^H \frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} \middle| \mathcal{F}_{k-1} \right] \\ &\leq 18H^2 \log\left(\frac{1}{\delta}\right) + 2HS(M+1)A + 4\sqrt{H^2S(M+1)AK} \\ &= O\left(H\left(SMA + H \log\left(\frac{1}{\delta}\right)\right) + \sqrt{H^2SMAK}\right) \\ &= \tilde{O}\left(\sqrt{H^2SMAK}\right) \end{aligned}$$

Proof. We start by rewriting the sum as follows:

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[\frac{z_{i,h}^k}{\sqrt{n_h^k(s_h^k, a_h^k, i) - 1}} \middle| \mathcal{F}_{k-1} \right] &= \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\sum_{i=1}^{M+1} z_{i,h}^k \frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, i) - 1}} \middle| \mathcal{F}_{k-1} \right] \\ &= \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\mathbb{E}_{x_h^k, z_{i,h}^k} \left[\frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} \right] \middle| \mathcal{F}_{k-1} \right] \\ &= \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} \middle| \mathcal{F}_{k-1} \right] \\ &= \sum_{k=1}^K \mathbb{E} \left[\sum_{h=1}^H \frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} \middle| \mathcal{F}_{k-1} \right], \end{aligned}$$

which proves the first equality. Now, defining $Y_k = \sum_{h=1}^H \frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}}$, which is \mathcal{F}_k -measurable and bounded almost surely in $[0, H]$, we can apply Lemma 27 of (Efroni et al., 2021) and get that for any $\delta > 0$, with probability at least $1 - \delta$,

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^{M+1} \mathbb{E} \left[\frac{z_{i,h}^k}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} \middle| \mathcal{F}_{k-1} \right] &\leq \left(1 + \frac{1}{2H}\right) \sum_{k=1}^K \sum_{h=1}^H \frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} + 2(2H+1)^2 \log \frac{1}{\delta} \\ &\leq 2 \sum_{k=1}^K \sum_{h=1}^H \frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} + 18H^2 \log \frac{1}{\delta}. \end{aligned}$$

Finally, observing that every time a context-state-action is visited, its count increases, we can bound the sum by

$$\begin{aligned}
 \sum_{k=1}^K \sum_{h=1}^H \frac{1}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} &= \sum_{k=1}^K \sum_{h=1}^H \sum_{x \in \mathcal{X}} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \frac{1_{\{x_h^k = x, s_h^k = s, a_h^k = a\}}}{\sqrt{n_h^k(s_h^k, a_h^k, x_h^k) - 1}} \\
 &= \sum_{h=1}^H \sum_{x \in \mathcal{X}} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \left(1 + \sum_{n=1}^{n_h^K(s, a, x)} \frac{1}{n} \right) \\
 &= HS(M+1)A + \sum_{h=1}^H \sum_{x \in \mathcal{X}} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} 2\sqrt{n_h^K(s, a, x)} \\
 &= HS(M+1)A + 2 \sqrt{HS(M+1)A \underbrace{\sum_{h=1}^H \sum_{x \in \mathcal{X}} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} n_h^K(s, a, x)}_{=HK}} \quad (\text{Cauchy Schwartz}) \\
 &= HS(M+1)A + 2\sqrt{H^2S(M+1)AK}.
 \end{aligned}$$

Substituting this bound concludes the proof. \square

Lemma E.8 (Elliptical potential lemma, Abbasi-Yadkori et al. (2011)). *Let $\{x_t, g_t\}_{t=1}^T$ be a sequence in \mathbb{R}^d such that $\|x_t\|_2 \leq L$ for all $t \leq T$ and let $\mathbf{V}_t = \lambda I + \sum_{s=1}^t x_s x_s^T$. Then,*

$$\sum_{t=1}^n \min\left\{ \|x_t\|_{\mathbf{V}_t^{-1}}^2, 1 \right\} \leq 2d \log \frac{\lambda d + nL^2}{\lambda d}$$

Corollary E.9. *Let $\{x_h^k, g_{k,1,h,2}^k\}_{k=1}^K$ be a sequence in \mathbb{R}^d such that $\|x_h^k\|_2 \leq L$ for all k, h and let $\mathbf{V}_k = \lambda I + \sum_{k^0=1}^k \sum_{h=1}^H x_h^{k^0} x_h^{k^0 T}$. Then,*

$$\sum_{k^0=1}^k \sum_{h=1}^H \left\| x_h^{k^0} \right\|_{\mathbf{V}_k^{-1}}^2 \leq \frac{\sqrt{2KH^2d \log \frac{\lambda d + kL^2}{\lambda d}}}{\max\{1, L/\lambda g\}}.$$

Proof. Define the matrices $\mathbf{V}_{k,h} = \sum_{k^0=1}^k x_h^{k^0} x_h^{k^0 T}$; clearly, it holds that $\mathbf{V}_k = \sum_{h=1}^H \mathbf{V}_{k,h}$ for all k, h , and thus, by applying Lemma E.8 for each of these matrices, we get

$$\sum_{k^0=1}^k \sum_{h=1}^H \min\left\{ \left\| x_h^{k^0} \right\|_{\mathbf{V}_k^{-1}}^2, 1 \right\} \leq \left(\sum_{k^0=1}^k \min\left\{ \left\| x_h^{k^0} \right\|_{\mathbf{V}_{k,h}^{-1}}^2, 1 \right\} \right) \sum_{h=1}^H 2d \log \frac{\lambda d + kL^2}{\lambda d} = 2dH \log \frac{\lambda d + kL^2}{\lambda d}.$$

Also, notice that if $\|x_h^k\|_2 \leq L$, then $\|x_h^k\|_{\mathbf{V}_k^{-1}}^2 \leq \frac{L^2}{\lambda_{\min}(\mathbf{V}_k)} \leq \frac{L^2}{\lambda}$, and thus

$$\left\| x_h^{k^0} \right\|_{\mathbf{V}_k^{-1}}^2 \leq \frac{\min\left\{ \left\| x_h^{k^0} \right\|_{\mathbf{V}_{k^0}^{-1}}^2, 1 \right\}}{\max\{1, L^2/\lambda g\}}.$$

Finally, the desired result is achieved by the Cauchy-Schwartz inequality:

$$\begin{aligned}
 \sum_{k^0=1}^k \sum_{h=1}^H \left\| x_h^{k^0} \right\|_{\mathbf{V}_k^{-1}}^2 &\leq \sqrt{KH \sum_{k^0=1}^k \sum_{h=1}^H \left\| x_h^{k^0} \right\|_{\mathbf{V}_{k^0}^{-1}}^2} \\
 &\leq \frac{\sqrt{KH \sum_{k^0=1}^k \sum_{h=1}^H \min\left\{ \left\| x_h^{k^0} \right\|_{\mathbf{V}_{k^0}^{-1}}^2, 1 \right\}}}{\max\{1, L/\lambda g\}} \\
 &\leq \frac{\sqrt{2KH^2d \log \frac{\lambda d + kL^2}{\lambda d}}}{\max\{1, L/\lambda g\}}
 \end{aligned}$$

Algorithm 4 Optimistic Threshold Planner for Logistic DCMDPs

```

1: require: Optimistic reward  $\bar{r}$ , estimated transition  $\hat{P}$ , and rectangular confidence set  $B_k(\delta)$  for  $\hat{f}_T$ .
2: init:  $\bar{V}_H(s, \mathbf{C}) = 0$ , for all  $s, \mathbf{C} \in \mathcal{S} \times \mathcal{I}^k$ 
3: for  $h = H - 1, \dots, 1$  do
4:   for each  $s \in \mathcal{S}, \mathbf{C}_h \in \{ \mathbf{C}(\sigma(\tau_h; \hat{f}_T)) : \tau_h \in H_h \}$  do
5:      $\mathbf{C}_{h+1} := \alpha \mathbf{C}_h + [l_h^k, u_h^k]$ 
6:      $\bar{Q}_i(s, a, \mathbf{C}_h) = \bar{r}_i(s, a) + E_{s^0 \sim \hat{P}_i(j_s, a)}[\bar{V}_{h+1}(s^0, \mathbf{C}_{h+1})]$  // State-action optimistic value
7:      $\bar{V}_h(s, \mathbf{C}_h) = \min \left\{ \max_{a \in \mathcal{A}, t \in \mathcal{T}} (\bar{Q}) \sum_{i=0}^M z_i(\mathbf{th}_t(\bar{Q}, \mathbf{C}_h)) \bar{Q}_i(s, a, \mathbf{C}_h), H \right\}$  // Lemma 5.6
8:      $\bar{\pi}(s, \mathbf{C}_h) \in \arg \max_{a \in \mathcal{A}} \max_{t \in \mathcal{T}} (\bar{Q}) \sum_{i=0}^M z_i(\mathbf{th}_t(\bar{Q}, \mathbf{C}_h)) \bar{Q}_i(s, a, \mathbf{C}_h)$ 
9:   end for
10: end for
11: Output  $\bar{\pi}(s, \tau) = \bar{\pi}(s, \mathbf{C}(\sigma(\tau)))$ 

```

□

F. Threshold Optimistic Planning

F.1. Proof of Threshold Optimism – Lemma 5.6

Lemma 5.6 (Threshold Optimism). *Let $\mathbf{Q} \in \mathcal{R}^{M+1}$. For any $\mathbf{x} \in \mathcal{R}^{M+1}$ such that $x_i = 0$ define $f(\mathbf{x}) = \sum_{i=1}^{M+1} z_i(\mathbf{x}) Q_i$. Let $\mathbf{C} = [l, u] \in \mathcal{R}^{M+1} \times \mathcal{R}^{M+1}$ and assume that $l < u$. Then, there exists $t \in \mathcal{T}(\mathbf{Q})$ such that $\mathbf{th}_t(\mathbf{Q}, \mathbf{C}) \in \arg \max_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x})$.*

Proof. For brevity, throughout the proof, we assume that $\eta = 1$, namely, $z_i(\mathbf{x}) = \frac{\exp(x_i)}{1 + \sum_{m=1}^M \exp(x_m)}$. This has no impact on the proof, since one can always denote $[\tilde{l}, \tilde{u}] = [\eta l, \eta u]$ and follow the rest of the proof with the modified intervals.

Let $X = \arg \max_{\mathbf{x} \in \mathcal{R}} f(\mathbf{x})$. We start by showing that there exists at least one solution at the extreme points of \mathcal{R} . We then show that solutions at the extreme points have a thresholding behavior.

Part 1. We first show that $X \setminus \text{ext}(\mathcal{R}) \neq \emptyset$, i.e., there exists $\mathbf{x} \in X$ that is an extreme point of the set \mathcal{R} . Note that f is continuous and \mathcal{R} is a compact set, therefore X is nonempty.

Let $\mathbf{x} \in X$ and choose some $k \in [M]$. We show that by replacing x_k by either l_k or u_k , we get another solution at X . Repeatedly doing so for all $k \in [M]$ will lead to $\mathbf{x} \in \text{ext}(\mathcal{R})$ and conclude this part of the proof.

We now fix $x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_M$ and study $f(\mathbf{x})$ only as a function of x_k . We also use the convention, $x_0 = 0$. Then,

$$\begin{aligned}
 f(\mathbf{x}) &= \sum_{i=0}^M z_i(\mathbf{x}) v_i \\
 &= \sum_{i=0}^M \frac{\exp(x_i)}{\sum_{j=0}^M \exp(x_j)} v_i \\
 &= \frac{\exp(x_k) v_k}{\sum_{j=0}^M \exp(x_j)} + \frac{\sum_{i \neq k} \exp(x_i) v_i}{\sum_{j=0}^M \exp(x_j)} \\
 &= \frac{\exp(x_k) v_k}{\sum_{j=0}^M \exp(x_j)} + \left(1 - \frac{\exp(x_k)}{\sum_{j=0}^M \exp(x_j)} \right) \frac{\sum_{i \neq k} \exp(x_i) v_i}{\sum_{j \neq k} \exp(x_j)}.
 \end{aligned}$$

Denote $\lambda(x_k) = \frac{\exp(x_k)}{\sum_{j=0}^M \exp(x_j)}$, and $v_{\text{ref}} = \frac{\sum_{i \neq k} \exp(x_i) v_i}{\sum_{j \neq k} \exp(x_j)}$. Then,

$$f(\mathbf{x}) = \lambda(x_k) v_k + (1 - \lambda(x_k)) v_{\text{ref}}.$$

Note that, since we fixed $x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_M$, then v_{ref} is constant (does not depend on x_k). Also, $\lambda(x_k)$ is a strictly monotonically increasing function in x_k and $f(\mathbf{x})$ is linear in $\lambda(x_k)$. Hence $\max_{x_k} f(\mathbf{x})$ is achieved either for $\tilde{x}_k = \arg \min_{x_k \in [l_k, u_k]} \lambda(x_k) = l_k$ or $\tilde{x}_k = \arg \max_{x_k \in [l_k, u_k]} \lambda(x_k) = u_k$. Denoting the solution that replaces x_k with the maximizer \tilde{x}_k by $\tilde{\mathbf{x}}$, we get that $f(\tilde{\mathbf{x}}) \geq f(\mathbf{x})$, but since $\mathbf{x} \in X$, so does $\tilde{\mathbf{x}} \in X$. Following this process for all $k \in [M]$ leads to an optimal $\tilde{\mathbf{x}} \in \text{ext}(\mathcal{R})$ and thus $X \setminus \text{ext}(\mathcal{R}) \neq \emptyset$.

Part 2. For the next part of the proof, we show that there exists an optimal solution that is a threshold function. Without loss of generality, assume that (v_1, \dots, v_M) are sorted in ascending order, such that $v_1 \leq v_2 \leq \dots \leq v_M$. Let $\mathbf{x} \in X \setminus \text{ext}(\mathcal{R})$, and assume by contradiction there exists $i, j \in [M], i < j$, such that $x_i = u_i, x_j = l_j$ and $v_i < v_j$. Denote

$$\begin{aligned} \epsilon_i &= \min\{x_i - \log(\exp(x_i) + \exp(x_j) - \exp(u_j)), x_i - l_i\} \\ \epsilon_j &= \log(\exp(x_i) + \exp(x_j) - \exp(x_i - \epsilon_i)) - x_j \end{aligned}$$

and let $\tilde{\mathbf{x}} = \mathbf{x} - \epsilon_i \mathbf{e}_i + \epsilon_j \mathbf{e}_j$. Then, ϵ_i, ϵ_j enjoy the following properties.

1. $\epsilon_i, \epsilon_j > 0$, since

$$\begin{aligned} \epsilon_i &= x_i - \log\left(\exp(x_i) + \underbrace{\exp(x_j) - \exp(u_j)}_{>0}\right) > 0, \\ \epsilon_j &= \log\left(\exp(x_i) + \underbrace{\exp(x_j) - \exp(x_i - \epsilon_i)}_{>0}\right) - x_j > 0. \end{aligned}$$

2. By definition, $\epsilon_i = x_i - l_i$ by definition, and $\epsilon_j = u_j - x_j$, since by substituting ϵ_i , we get

$$\begin{aligned} \epsilon_j &= \log(\exp(x_i) + \exp(x_j) - \exp(x_i - [x_i - \log(\exp(x_i) + \exp(x_j) - \exp(u_j))])) - x_j \\ &= \log(\exp(x_i) + \exp(x_j) - [\exp(x_i) + \exp(x_j) - \exp(u_j)]) - x_j \\ &= u_j - x_j. \end{aligned}$$

In particular, given that $\epsilon_i, \epsilon_j > 0$, it implies that $l_i = x_i - \epsilon_i = u_i$ and $l_j = x_j + \epsilon_j = u_j$.

3. The total weight of i, j is preserved

$$\begin{aligned} \exp(x_i - \epsilon_i) + \exp(x_j + \epsilon_j) &= \exp(x_i - \epsilon_i) + [\exp(x_i) + \exp(x_j) - \exp(x_i - \epsilon_i)] \\ &= \exp(x_i) + \exp(x_j) \end{aligned}$$

Given these properties, $\tilde{\mathbf{x}}$ is a valid solution for which we have that

$$\begin{aligned} f(\tilde{\mathbf{x}}) &= \frac{\exp(\tilde{x}_i) v_i}{1 + \sum_{k=0}^M \exp(\tilde{x}_k)} + \frac{\exp(\tilde{x}_j) v_j}{1 + \sum_{k=0}^M \exp(\tilde{x}_k)} + \frac{\sum_{k \neq i, j} \exp(\tilde{x}_k) v_k}{1 + \sum_{k=1}^M \exp(\tilde{x}_k)} \\ &= \frac{\exp(\tilde{x}_i) v_i}{1 + \sum_{k=0}^M \exp(x_k)} + \frac{\exp(\tilde{x}_j) v_j}{1 + \sum_{k=0}^M \exp(x_k)} + \frac{\sum_{k \neq i, j} \exp(x_k) v_k}{1 + \sum_{k=1}^M \exp(x_k)}. \end{aligned} \quad (\text{By property (3)})$$

Therefore,

$$f(\tilde{\mathbf{x}}) - f(\mathbf{x}) = \frac{\exp(x_i - \epsilon_i) v_i + \exp(x_j + \epsilon_j) v_j - \exp(x_i) v_i - \exp(x_j) v_j}{1 + \sum_{k=0}^M \exp(x_k)}.$$

Considering the numerator, we have that

$$\begin{aligned}
 & \exp(x_i - \epsilon_i)v_i + \exp(x_j + \epsilon_j)v_j - \exp(x_i)v_i - \exp(x_j)v_j \\
 &= \exp(x_i - \epsilon_i)v_i + (\exp(x_i) + \exp(x_j) - \exp(x_i - \epsilon_i))v_j - \exp(x_i)v_i - \exp(x_j)v_j \quad (\text{By property (3)}) \\
 &= \exp(x_i - \epsilon_i)v_i + \exp(x_i)v_j - \exp(x_i - \epsilon_i)v_j - \exp(x_i)v_i \\
 &= (\exp(x_i) - \exp(x_i - \epsilon_i))(v_j - v_i) \\
 &> 0,
 \end{aligned}$$

where the inequality is since $\epsilon_i > 0$ and $v_i < v_j$. That is, $f(\tilde{\mathbf{x}}) > f(\mathbf{x})$, in contradiction to $\mathbf{x} \succeq X$. To summarize, we prove that for any $\mathbf{x} \succeq X \setminus \text{ext}(\mathcal{R})$, if $v_i < v_j$, then $x_i < x_j$, which corresponds to a thresholding function. All that is left is to prove that if $v_i = v_j = v$, there exists a solution $\mathbf{x} \succeq X \setminus \text{ext}(\mathcal{R})$ such that either $x_i = u_i, x_j = u_j$ or $x_i = l_i, x_j = l_j$. To show this, we follow a similar path to the first part of the proof and write

$$\begin{aligned}
 f(\mathbf{x}) &= \sum_{k=0}^M z_i(\mathbf{x})v_i \\
 &= \frac{\exp(x_i)v_i + \exp(x_j)v_j}{\sum_{k=0}^M \exp(x_k)} + \left(1 - \frac{\exp(x_i) + \exp(x_j)}{\sum_{k=0}^M \exp(x_k)}\right) \frac{\sum_{k \notin \{i,j\}} \exp(x_k)v_k}{\sum_{k \notin \{i,j\}} \exp(x_k)} \\
 &= \frac{\exp(x_i) + \exp(x_j)}{\sum_{k=0}^M \exp(x_k)}v + \left(1 - \frac{\exp(x_i) + \exp(x_j)}{\sum_{k=0}^M \exp(x_k)}\right) \frac{\sum_{k \notin \{i,j\}} \exp(x_k)v_k}{\sum_{k \notin \{i,j\}} \exp(x_k)} \quad (v_i = v_j = v)
 \end{aligned}$$

Now, denoting $\lambda(x_i, x_j) = \frac{\exp(x_i) + \exp(x_j)}{\sum_{k=0}^M \exp(x_k)}$, and $v_{\text{ref}} = \frac{\sum_{k \notin \{i,j\}} \exp(x_k)v_k}{\sum_{k \notin \{i,j\}} \exp(x_k)}$, we can follow the exact same line of the proof as the first part, and conclude that there exist another solution $\tilde{\mathbf{x}}$ in which $\lambda(x_i, x_j)$ is either maximized or minimized – either $x_i = u_i, x_j = u_j$ or $x_i = l_i, x_j = l_j$.

This completes the proof. □

G. Confidence Sets

We require the following quantities, used by Amani and Thrampoulidis (2021) or adapted from Abeille et al. (2021). First, recall that for any $\mathbf{f} \in \mathbb{R}^{M(M+1)SAH}$ and $\mathbf{d} \in \mathbb{R}^{MSAH}$, we have

$$\mathbf{A}(\mathbf{d}, \mathbf{f}) := \text{diag}(\mathbf{z}(\mathbf{d}, \mathbf{f})) \quad \mathbf{z}(\mathbf{d}, \mathbf{f})\mathbf{z}(\mathbf{d}, \mathbf{f})^T$$

Also, for any $\mathbf{f} \in \mathbb{R}^{M(M+1)SAH}$, define

$$\mathbf{g}_k(\mathbf{f}) := \lambda \mathbf{f} + \sum_{k^0=1}^{k-1} \sum_{h=1}^H z_i(\mathbf{d}_h^{k^0}, \mathbf{f}) \quad \mathbf{d}_h^{k^0}, \quad \text{and} \quad \mathbf{H}_k(\mathbf{f}) := \lambda \mathbf{I} + \sum_{k^0=1}^{k-1} \sum_{h=1}^H \mathbf{A}(\mathbf{d}_h^{k^0}, \mathbf{f}) \quad \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T}.$$

Note that by definition

$$\Gamma_{\mathbf{f}} L_{\lambda}^k(\mathbf{f}) = \sum_{k^0=1}^{k-1} \sum_{h=1}^H \mathbf{m}_h^{k^0} \quad \mathbf{d}_h^{k^0} \quad \mathbf{g}_k(\mathbf{f}) \quad \text{and} \quad \Gamma_{\mathbf{f}}^2 L_{\lambda}^k(\mathbf{f}) = \mathbf{H}_k(\mathbf{f}) \quad (11)$$

Confidence Set

$$\mathcal{C}_k(\delta) := \left\{ \mathbf{f} \in F : \left\| \mathbf{g}_k(\mathbf{f}) - \mathbf{g}_k(\hat{\mathbf{f}}_t) \right\|_{\mathbf{H}_k^{-1}(\mathbf{f})} \leq \beta_k(\delta) \right\}, \quad (12)$$

where $\beta_k(\delta) = \frac{M^{3/2}(M+1)SAH}{\beta_{\lambda}} \left(\log \left(1 + \frac{k}{(M+1)SA\lambda} \right) + 2 \log \left(\frac{2}{\delta} \right) \right) + \sqrt{\frac{\lambda}{4M}} + \frac{1}{\lambda} L$.

Other Notations For any $\mathbf{f}_1, \mathbf{f}_2 \in \mathbb{R}^{M(M+1)SAH}$ and $\mathbf{d} \in \mathbb{R}^{MSAH}$, define

$$\begin{aligned} \mathbf{B}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) &:= \int_0^1 \mathbf{A}(\mathbf{d}, v\mathbf{f}_1 + (1-v)\mathbf{f}_2) dv, \\ \tilde{\mathbf{B}}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) &:= \int_0^1 (1-v)\mathbf{A}(\mathbf{d}, v\mathbf{f}_1 + (1-v)\mathbf{f}_2) dv, \\ \mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2) &:= \lambda \mathbf{I} + \sum_{k^0=1}^{k-1} \sum_{h=1}^H \mathbf{B}(\mathbf{d}_h^{k^0}, \mathbf{f}_1, \mathbf{f}_2) \quad \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T}, \\ \tilde{\mathbf{G}}_k(\mathbf{f}_1, \mathbf{f}_2) &:= \lambda \mathbf{I} + \sum_{k^0=1}^{k-1} \sum_{h=1}^H \tilde{\mathbf{B}}(\mathbf{d}_h^{k^0}, \mathbf{f}_1, \mathbf{f}_2) \quad \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T}, \\ \mathbf{V}_k &:= \lambda \mathbf{I}_{(M+1)SAH} + \sum_{k^0=1}^{k-1} \sum_{h=1}^H \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T}. \end{aligned}$$

Note that we stray from the notation of \mathbf{V}_k in (Amani and Thrampoulidis, 2021), by removing the factor of κ from the regularization term.

G.1. Useful Lemmas

We now provide a list of lemmas required for providing confidence intervals for the logistic history dependent transition model in Appendices C and D. In what follows we will use the following expression:

$$d_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) := \left\| (\mathbf{f}_1 \quad \mathbf{f}_2)^T \mathbf{d} \right\|_2 \leq \kappa \mathbf{f}_1 \quad \mathbf{f}_2 \kappa \mathbf{d} \kappa_2 \leq L, \quad (13)$$

To this end, the following properties hold:

Lemma G.1 (Amani and Thrampoulidis (2021), Lemma 2). *For any $\mathbf{f}_1, \mathbf{f}_2 \in \mathbb{R}^{M(M+1)SAH}$ and $\mathbf{d} \in \mathbb{R}^{MSAH}$*

$$\mathbf{z}(\mathbf{d}, \mathbf{f}_1) - \mathbf{z}(\mathbf{d}, \mathbf{f}_2) = [\mathbf{B}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) \quad \mathbf{d}](\mathbf{f}_1 \quad \mathbf{f}_2)$$

Lemma G.2 (Amani and Thrampoulidis (2021), Lemma 3). *For any $\mathbf{f}_1, \mathbf{f}_2 \in \mathbb{R}^{M(M+1)SAH}$,*

$$\mathbf{g}_k(\mathbf{f}_1) - \mathbf{g}_k(\mathbf{f}_2) = \mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2)(\mathbf{f}_1 - \mathbf{f}_2)$$

Lemma G.3 (Amani and Thrampoulidis (2021), Lemma 4). *For any $\mathbf{f}_1, \mathbf{f}_2 \in F$, it holds that $(1 + 2L)^{-1} \mathbf{H}_k(\mathbf{f}_1) - \mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2)$ and $(1 + 2L)^{-1} \mathbf{H}_k(\mathbf{f}_2) - \mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2)$.*

Lemma G.4 (Amani and Thrampoulidis (2021), Lemma 5). *For any $\mathbf{f} \in \mathbb{R}^{M(M+1)SAH}$ and $\mathbf{d} \in \mathbb{R}^{MSAH}$, the matrix $\mathbf{A}(\mathbf{d}, \mathbf{f})$ is strictly diagonally dominant and thus positive definite.*

Lemma G.5 (Amani and Thrampoulidis (2021), Theorem 1). *Let $\delta \in (0, 1)$. With probability at least $1 - \delta$, for all $k \geq 1$, it holds that $\mathbf{f} \in \mathcal{C}_k(\delta)$.*

Remark G.6. Notice that Lemma G.4 also implies that all matrices $\mathbf{B}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2), \tilde{\mathbf{B}}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2), \mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2), \tilde{\mathbf{G}}_k(\mathbf{f}_1, \mathbf{f}_2)$ are positive definite. Moreover, it implies that $\mathbf{B}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) \succeq \tilde{\mathbf{B}}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2)$ (since $1 - v \in [0, 1]$), and therefore, $\mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2) \succeq \tilde{\mathbf{G}}_k(\mathbf{f}_1, \mathbf{f}_2)$.

Relying on our different definition for V_k , the next lemma allows us to gain dependence on the derivative at the *real* latent features \mathbf{f} , instead of the worst-case derivative as in (Amani and Thrampoulidis, 2021).

Lemma G.7 (Connection between Local and Global Design Matrices). *For any $\mathbf{f} \in F$, denote $\kappa(\mathbf{f}) = \frac{1}{\inf_{\mathbf{d} \in \mathcal{D}} \lambda_{\min} \mathbf{f} \mathbf{A}(\mathbf{d}; \mathbf{f}) \mathbf{g}}$. It holds that*

$$\kappa(\mathbf{f}) \mathbf{H}_k(\mathbf{f}) \succeq \mathbf{I}_M - \mathbf{V}_k,$$

and specifically for $\mathbf{f} \in F$,

$$\kappa \mathbf{H}_k(\mathbf{f}) \succeq \mathbf{I}_M - \mathbf{V}_k.$$

Proof. For clarity, we state the dimension of the identity matrices throughout the proof. Throughout the analysis, recall that if $\mathbf{A}, \mathbf{B}, \mathbf{C} \succeq 0$ and $\mathbf{A} \succeq \mathbf{B}$ then $\mathbf{A} + \mathbf{C} \succeq \mathbf{B} + \mathbf{C}$. For any $\mathbf{f} \in F$,

$$\begin{aligned} \mathbf{H}_k(\mathbf{f}) &= \lambda \mathbf{I}_{M(M+1)SAH} + \sum_{k^0=1}^k \sum_{h=1}^H \mathbf{A}(\mathbf{d}_h^{k^0}, \mathbf{f}) - \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T} \\ &\succeq \lambda \mathbf{I}_{M(M+1)SAH} + \sum_{k^0=1}^k \sum_{h=1}^H \lambda_{\min} \{ \mathbf{A}(\mathbf{d}_h^{k^0}; \mathbf{f}) \} \mathbf{I}_{M+1} - \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T} \\ &\succeq \lambda \mathbf{I}_{M(M+1)SAH} + \left(\inf_{\mathbf{d} \in \mathcal{D}} \lambda_{\min} \mathbf{f} \mathbf{A}(\mathbf{d}; \mathbf{f}) \mathbf{g} \right) \sum_{k^0=1}^k \sum_{h=1}^H \mathbf{I}_{M+1} - \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T} \\ &= \lambda \mathbf{I}_{M(M+1)SAH} + \frac{1}{\kappa(\mathbf{f})} \sum_{k^0=1}^k \sum_{h=1}^H \mathbf{I}_{M+1} - \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T} \\ &= \frac{1}{\kappa(\mathbf{f})} \mathbf{I}_M \left(\kappa(\mathbf{f}) \lambda \mathbf{I}_{(M+1)SAH} + \sum_{k^0=1}^k \sum_{h=1}^H \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T} \right) \\ &\stackrel{(\cdot)}{=} \frac{1}{\kappa(\mathbf{f})} \mathbf{I}_M \left(\lambda \mathbf{I}_{(M+1)SAH} + \sum_{k^0=1}^k \sum_{h=1}^H \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T} \right) \\ &= \frac{1}{\kappa(\mathbf{f})} \mathbf{I}_M - \mathbf{V}_k, \end{aligned}$$

where (\cdot) holds since $\kappa(\mathbf{f}) \geq 1$ (see, e.g., eq. (29) of Amani and Thrampoulidis (2021) when fixing the set of possible parameters F to be a singleton $F = \mathbf{f} \mathbf{f} \mathbf{g}$).

Finally, we conclude the proof by noting that for $\mathbf{f} \in F$, it holds by definition that $\kappa = \kappa(\mathbf{f})$. \square

Lemma G.8. *For all $\mathbf{d} \in \mathbb{R}^{MSAH}$ such that $k \|\mathbf{d}\|_2 \leq 1$, $k \geq 1$ and $\mathbf{f} \in \mathcal{C}_k(\delta)$, if $\mathbf{f} \in \mathcal{C}_k(\delta)$ then*

$$\|\kappa \mathbf{z}(\mathbf{d}, \mathbf{f}) - \mathbf{z}(\mathbf{d}, \mathbf{f})\|_2 \leq 2\beta_k(\delta)(1 + 2L)^{\rho_-} \frac{1}{\kappa} k \|\mathbf{d}\|_{\mathbf{V}_k^{-1}}$$

Proof. Here, we closely follow the proof of Lemma 1 in (Amani and Thrapoulidis, 2021), with the exception that we apply Lemma G.7 to achieve dependence on κ . Specifically, we let $\kappa_{\max} := \sup_{\mathbf{d}, \mathbf{f} \in \mathcal{D}, \mathcal{F}} \lambda_{\max} \bar{f} \mathbf{A}(\mathbf{d}; \mathbf{f}) \mathbf{g}$. Notice that following (Amani and Thrapoulidis, 2021, Section 3), it holds that $\kappa_{\max} \geq 1$.

$$\begin{aligned}
 & \kappa z(\mathbf{d}, \mathbf{f}) - z(\mathbf{d}, \mathbf{f}) \kappa_2 \\
 &= \kappa [\mathbf{B}(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d}] (\mathbf{f} - \mathbf{f}) \kappa \tag{Lemma G.1} \\
 &= \left\| [\mathbf{B}(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d}] \mathbf{G}_k^{-1/2}(\mathbf{f}, \mathbf{f}) \mathbf{G}_k^{1/2}(\mathbf{f}, \mathbf{f}) (\mathbf{f} - \mathbf{f}) \right\| \\
 & \quad \left\| [\mathbf{B}(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d}] \mathbf{G}_k^{-1/2}(\mathbf{f}, \mathbf{f}) \right\| \kappa \mathbf{f} - \mathbf{f} \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \tag{Cauchy-Schwartz} \\
 &= \left\| [\mathbf{B}(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d}] \mathbf{G}_k^{-1/2}(\mathbf{f}, \mathbf{f}) \right\| \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \tag{Lemma G.2} \\
 &= \sqrt{\lambda_{\max}([\mathbf{B}(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d}] \mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f}) [\mathbf{B}(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d}]) \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})}} \\
 &= \sqrt{\lambda_{\max}(\mathbf{G}_k^{-1/2}(\mathbf{f}, \mathbf{f}) [\mathbf{B}^T(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d}] [\mathbf{B}(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d}]^T \mathbf{G}_k^{-1/2}(\mathbf{f}, \mathbf{f})) \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})}} \\
 & \quad \text{(cyclic property of } \lambda_{\max}) \\
 &= \sqrt{\lambda_{\max}(\mathbf{G}_k^{-1/2}(\mathbf{f}, \mathbf{f}) [\mathbf{B}^T(\mathbf{d}, \mathbf{f}, \mathbf{f}) \mathbf{B}(\mathbf{d}, \mathbf{f}, \mathbf{f}) - \mathbf{d} \mathbf{d}^T] \mathbf{G}_k^{-1/2}(\mathbf{f}, \mathbf{f})) \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})}} \\
 & \quad \text{(mixed-product property)} \\
 &= \kappa_{\max} \sqrt{\lambda_{\max}(\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f}) [\mathbf{I}_M - \mathbf{d} \mathbf{d}^T])} \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \tag{definition of } \kappa_{\max}) \\
 &= \sqrt{\lambda_{\max}(\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f}) [\mathbf{I}_M - \mathbf{d}] [\mathbf{I}_M - \mathbf{d}^T])} \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \tag{cyclic property of } \lambda_{\max}, \text{ and } \kappa_{\max} \geq 1) \\
 &= \sqrt{\lambda_{\max}([\mathbf{I}_M - \mathbf{d}] \mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f}) [\mathbf{I}_M - \mathbf{d}^T])} \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \tag{mixed-product property} \\
 & \quad \rho_{1+2L} \sqrt{\lambda_{\max}([\mathbf{I}_M - \mathbf{d}] \mathbf{H}_k^{-1}(\mathbf{f}) [\mathbf{I}_M - \mathbf{d}^T])} \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \tag{Lemma G.3} \\
 & \quad \sqrt{\kappa(1+2L)} \sqrt{\lambda_{\max}([\mathbf{I}_M - \mathbf{d}] [\mathbf{I}_M - \mathbf{V}_k^{-1}(\mathbf{f})] [\mathbf{I}_M - \mathbf{d}^T])} \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \tag{Lemma G.7} \\
 &= \sqrt{\kappa(1+2L)} \sqrt{\lambda_{\max}(\mathbf{I}_M - \kappa \mathbf{d} \mathbf{d}^T \mathbf{V}_k^{-1})} \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \tag{mixed-product property} \\
 &= \sqrt{\kappa(1+2L)} \kappa \mathbf{d} \mathbf{d}^T \mathbf{V}_k^{-1} \kappa g_k(\mathbf{f}) - g_k(\mathbf{f}) \kappa_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \\
 &= \sqrt{\kappa(1+2L)} \kappa \mathbf{d} \mathbf{d}^T \mathbf{V}_k^{-1} \left\| g_k(\mathbf{f}) - g_k(\hat{\mathbf{f}}) + g_k(\hat{\mathbf{f}}) - g_k(\mathbf{f}) \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \\
 &= \sqrt{\kappa(1+2L)} \kappa \mathbf{d} \mathbf{d}^T \mathbf{V}_k^{-1} \left[\left\| g_k(\mathbf{f}) - g_k(\hat{\mathbf{f}}) \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} + \left\| g_k(\hat{\mathbf{f}}) - g_k(\mathbf{f}) \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \mathbf{f})} \right] \\
 & \quad (1+2L) \rho_{\kappa}^{-} \kappa \mathbf{d} \mathbf{d}^T \mathbf{V}_k^{-1} \left[\left\| g_k(\mathbf{f}) - g_k(\hat{\mathbf{f}}) \right\|_{\mathbf{H}_k^{-1}(\mathbf{f})} + \left\| g_k(\hat{\mathbf{f}}) - g_k(\mathbf{f}) \right\|_{\mathbf{H}_k^{-1}(\mathbf{f})} \right] \tag{Lemma G.3} \\
 & \quad 2\beta_k(\delta)(1+2L) \rho_{\kappa}^{-} \kappa \mathbf{d} \mathbf{d}^T \mathbf{V}_k^{-1}. \tag{f, f \in \mathcal{C}_k(\delta)}
 \end{aligned}$$

The *cyclic property* of λ_{\max} refers to the fact that for two matrices $\mathbf{M}_1, \mathbf{M}_2$, the eigenvalues of $\mathbf{M}_1 \mathbf{M}_2$ are the same as the eigenvalues of $\mathbf{M}_2 \mathbf{M}_1$, and thus the same hold for the maximal eigenvalue. \square

Lemma G.9 (Adaptation of (Abeille et al., 2021), Lemma 8, to the multinomial case in Amani and Thrapoulidis (2021), Lemma 13). *For any $\mathbf{f}_1, \mathbf{f}_2 \in \mathcal{F}$, it holds that $(2+2L)^{-1} \mathbf{H}_k(\mathbf{f}_1) - \tilde{\mathbf{G}}_k(\mathbf{f}_1, \mathbf{f}_2)$ and $(2+2L)^{-1} \mathbf{H}_k(\mathbf{f}_2) - \tilde{\mathbf{G}}_k(\mathbf{f}_1, \mathbf{f}_2)$.*

Proof. According to Sun and Tran-Dinh (2019)[Eq. 16], for any $\mathbf{d} \in \mathcal{R}^{M_{SAH}}$, $\mathbf{f}_1, \mathbf{f}_2 \in \mathcal{R}^{M(M+1)_{SAH}}$, and for any $v \in [0, 1]$, we have that $r^{-2} f(vx + (1-v)y) \leq e^{-vd_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2)} r^{-2} f(y)$, where $d_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2)$ is defined in Equation (13).

Thus,

$$\int_0^1 (1-v)r^2 f(vx + (1-v)y)dv = r^2 f(y) \int_0^1 (1-v)e^{-vd_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2)} dv,$$

and replacing with the notation $r^2 f(\mathbf{x}) = \mathbf{A}(\mathbf{d}, \mathbf{x})$, we get

$$\tilde{\mathbf{B}}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) = \int_0^1 (1-v)\mathbf{A}(\mathbf{d}, v\mathbf{f}_1 + (1-v)\mathbf{f}_2)dv = \mathbf{A}(\mathbf{d}, \mathbf{f}_2) \int_0^1 (1-v)e^{-vd_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2)} dv,$$

Integrating the RHS by parts,

$$\tilde{\mathbf{B}}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) = \left(\frac{1}{d_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2)} + \frac{e^{-d_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2)}}{(d_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2))^2} \right) r^2 f(y) = g(d_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2))\mathbf{A}(\mathbf{d}, \mathbf{f}_2),$$

where we defined $g(z) = \frac{1}{z} \left(1 + \frac{e^{-z}}{z} \right)$.

Next, by Abeille et al. (2021)[Lemma 10], for all $z > 0$, it holds that $g(z) \geq \frac{1}{2+z}$, and therefore,

$$\tilde{\mathbf{B}}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) \geq (2 + d_2(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2))^{-1} \mathbf{A}(\mathbf{d}, \mathbf{f}_2) \geq (2 + 2L)^{-1} \mathbf{A}(\mathbf{d}, \mathbf{f}_2), \quad (14)$$

where the last inequality follows is due to Equation (13).

Plugging in Equation (14) with the definition of $\tilde{\mathbf{G}}_k(\mathbf{f}_1, \mathbf{f}_2)$,

$$\begin{aligned} \tilde{\mathbf{G}}_k(\mathbf{f}_1, \mathbf{f}_2) &= \lambda \mathbf{I} + \sum_{k^0=1}^k \sum_{h=1}^H \tilde{\mathbf{B}}(\mathbf{d}_h^{k^0}, \mathbf{f}_1, \mathbf{f}_2) \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T} \\ &\geq (2 + 2L)^{-1} \left(\lambda \mathbf{I} + \sum_{k^0=1}^k \sum_{h=1}^H \mathbf{A}(\mathbf{d}_h^{k^0}, \mathbf{f}_2) \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0 T} \right) \\ &= (2 + 2L)^{-1} \mathbf{H}_k(\mathbf{f}_2). \end{aligned}$$

By the symmetry in $\mathbf{f}_1, \mathbf{f}_2$ in the definition of $\tilde{\mathbf{B}}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2)$, we can similarly prove that $(2 + 2L)^{-1} \mathbf{H}_k(\mathbf{f}_1) \leq \tilde{\mathbf{G}}_k(\mathbf{f}_1, \mathbf{f}_2)$ \square

G.2. Convex Relaxation

Similar to Abeille et al. (2021), we define the convex relaxation of the set $C_k(\delta)$ by

$$E_k(\delta) = \left\{ \mathbf{f} \geq F : L_\lambda^k(\hat{\mathbf{f}}_k) \leq L_\lambda^k(\mathbf{f}) \leq \xi^2(\delta) \right\} \quad \text{where } \xi(\delta) = \beta_k(\delta) + \sqrt{\frac{HM}{\lambda}} \beta_k^2(\delta). \quad (15)$$

The next proposition is an adaptation of (Abeille et al., 2021)[Lemma 1] to the multinomial setting of (Amani and Thrampoulidis, 2021). Importantly, this proposition provides a confidence interval for the convex relaxation set Equation (15), which serves as the basis for the tractable estimator in Section 5.1.

Proposition G.10. *Let $\delta \geq (0, 1)$.*

1. $C_k(\delta) \subseteq E_k(\delta)$ for all $k \geq 1$ and therefore, w.p. $1 - \delta$, $\mathbf{f} \geq E_k(\delta)$ for all $k \geq 1$.
2. With probability $1 - \delta$, it holds that

$$\exists \mathbf{f} \geq E_k(\delta), \quad k \mathbf{f} \leq \mathbf{f} \leq \mathbf{K}_{\mathbf{H}_k(\mathbf{f})} \left((2 + 2L)\beta_k(\delta) + \sqrt{2(1+L)}\xi_k(\delta) \right).$$

In particular, for $\bar{\mathbf{f}} \geq \arg \max_{\mathbf{f} \geq F} L_\lambda^k(\mathbf{f})$, with probability $1 - \delta$,

$$\|\bar{\mathbf{f}} - \mathbf{f}\|_{\mathbf{H}_k(\mathbf{f})} \leq (2 + 2L)\beta_k(\delta) + \sqrt{2(1+L)}\xi_k(\delta) + \gamma_k(\delta),$$

where $\gamma_k(\delta) := \left(2 + 2L + \sqrt{2(1+L)} \right) \beta_k(\delta) + \sqrt{\frac{2(1+L)HM}{\lambda}} \beta_k^2(\delta)$.

As in (Abeille et al., 2021), in order to prove Proposition G.10, we first require the following side-lemma:

Lemma G.11 (Counterpart of Abeille et al. (2021), Lemma 2). *Let $\delta \in (0, 1)$. For all $\mathbf{f} \in \mathcal{C}_k(\delta)$, it holds that*

$$\|\mathbf{g}_k(\mathbf{f}) - \mathbf{g}_k(\hat{\mathbf{f}}_k)\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \hat{\mathbf{f}}_k)} \leq \xi_k(\delta).$$

Proof. First notice that by Remark G.6, the norm w.r.t. the inverse is well-defined and we further have $\lambda_{\min}(\mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2)) \geq \lambda$. Next, we utilise Amani and Thrampoulidis (2021), (eq. 61), which states that for any $\mathbf{f}_1, \mathbf{f}_2 \in F$, $\mathbf{d} \in D$ (Recall that \mathbf{B} is symmetric between $\mathbf{f}_1, \mathbf{f}_2$):

$$\begin{aligned} \mathbf{B}(\mathbf{d}, \mathbf{f}_1, \mathbf{f}_2) &= \left(1 + \left\| \begin{bmatrix} \mathbf{f}_{11}^T \mathbf{d} & \mathbf{f}_{21}^T \mathbf{d} & \dots & \mathbf{f}_{1M}^T \mathbf{d} & \mathbf{f}_{2M}^T \mathbf{d} \end{bmatrix} \right\| \right)^{-1} \mathbf{A}(\mathbf{d}, \mathbf{f}_1) && \text{(Amani and Thrampoulidis (2021), (eq. 61))} \\ &= \left(1 + \mathbf{1}_M^T \mathbf{d} \mathbf{K}_{\mathbf{G}_k^{-1}(\mathbf{f}_1, \mathbf{f}_2)} \mathbf{K} \begin{bmatrix} \mathbf{f}_1 & \mathbf{f}_2 \end{bmatrix} \mathbf{K}_{\mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2)} \right)^{-1} \mathbf{A}(\mathbf{d}, \mathbf{f}_1) && \text{(Cauchy-Schwartz)} \\ &= \left(1 + \sqrt{\frac{HM}{\lambda}} \mathbf{K} \begin{bmatrix} \mathbf{f}_1 & \mathbf{f}_2 \end{bmatrix} \mathbf{K}_{\mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2)} \right)^{-1} \mathbf{A}(\mathbf{d}, \mathbf{f}_1), && \text{(\mathbf{G}_k(\mathbf{f}_1, \mathbf{f}_2) \geq \lambda \mathbf{I})} \end{aligned}$$

where $j \in [1, \dots, M]$, \mathbf{f}_{ij} is the j -th coordinate of \mathbf{f}_i and $\mathbf{1}_M \in \mathbb{R}^M$ is a vector of ones.

Thus, we can write for any $\mathbf{f} \in \mathcal{C}_k(\delta)$

$$\begin{aligned} \mathbf{G}_k(\mathbf{f}, \hat{\mathbf{f}}_k) &= \lambda \mathbf{I} + \sum_{k^\circ=1}^k \sum_{h=1}^H \mathbf{B}(\mathbf{d}_h^{k^\circ}, \mathbf{f}, \hat{\mathbf{f}}_k) \mathbf{d}_h^{k^\circ} \mathbf{d}_h^{k^\circ T} \\ &= \lambda \mathbf{I} + \left(1 + \sqrt{\frac{HM}{\lambda}} \left\| \begin{bmatrix} \mathbf{f} & \hat{\mathbf{f}}_k \end{bmatrix} \right\|_{\mathbf{G}_k(\mathbf{f}, \hat{\mathbf{f}}_k)} \right)^{-1} \sum_{k^\circ=1}^k \sum_{h=1}^H \mathbf{A}(\mathbf{d}_h^{k^\circ}, \mathbf{f}) \mathbf{d}_h^{k^\circ} \mathbf{d}_h^{k^\circ T} \quad (\delta \mathbf{A}, \mathbf{B} \geq 0) \quad \mathbf{A} \geq \mathbf{B} \geq 0) \\ &= \left(1 + \sqrt{\frac{HM}{\lambda}} \left\| \begin{bmatrix} \mathbf{f} & \hat{\mathbf{f}}_k \end{bmatrix} \right\|_{\mathbf{G}_k(\mathbf{f}, \hat{\mathbf{f}}_k)} \right)^{-1} \left(\lambda \mathbf{I} + \sum_{k^\circ=1}^k \sum_{h=1}^H \mathbf{A}(\mathbf{d}_h^{k^\circ}, \mathbf{f}) \mathbf{d}_h^{k^\circ} \mathbf{d}_h^{k^\circ T} \right) \\ &= \left(1 + \sqrt{\frac{HM}{\lambda}} \left\| \begin{bmatrix} \mathbf{f} & \hat{\mathbf{f}}_k \end{bmatrix} \right\|_{\mathbf{G}_k(\mathbf{f}, \hat{\mathbf{f}}_k)} \right)^{-1} \mathbf{H}_k(\mathbf{f}) \\ &= \left(1 + \sqrt{\frac{HM}{\lambda}} \left\| \begin{bmatrix} \mathbf{g}_k(\mathbf{f}) & \mathbf{g}_k(\hat{\mathbf{f}}_k) \end{bmatrix} \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \hat{\mathbf{f}}_k)} \right)^{-1} \mathbf{H}_k(\mathbf{f}), && \text{(Lemma G.2)} \end{aligned}$$

Using this inequality, we get

$$\begin{aligned} \left\| \begin{bmatrix} \mathbf{g}_k(\mathbf{f}) & \mathbf{g}_k(\hat{\mathbf{f}}_k) \end{bmatrix} \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \hat{\mathbf{f}}_k)}^2 &= \left(1 + \sqrt{\frac{HM}{\lambda}} \left\| \begin{bmatrix} \mathbf{g}_k(\mathbf{f}) & \mathbf{g}_k(\hat{\mathbf{f}}_k) \end{bmatrix} \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \hat{\mathbf{f}}_k)} \right)^2 \left\| \begin{bmatrix} \mathbf{g}_k(\mathbf{f}) & \mathbf{g}_k(\hat{\mathbf{f}}_k) \end{bmatrix} \right\|_{\mathbf{H}_k^{-1}(\mathbf{f})}^2 \\ &= \left(1 + \sqrt{\frac{HM}{\lambda}} \left\| \begin{bmatrix} \mathbf{g}_k(\mathbf{f}) & \mathbf{g}_k(\hat{\mathbf{f}}_k) \end{bmatrix} \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \hat{\mathbf{f}}_k)} \right)^2 \beta_k^2(\delta). \quad (\mathbf{f} \in \mathcal{C}_k(\delta), \text{ see Equation (12)}) \end{aligned}$$

Solving this inequality finally yields the desired result (see, e.g. Abeille et al., 2021, Proposition 7):

$$\left\| \begin{bmatrix} \mathbf{g}_k(\mathbf{f}) & \mathbf{g}_k(\hat{\mathbf{f}}_k) \end{bmatrix} \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \hat{\mathbf{f}}_k)} \leq \beta_k(\delta) + \sqrt{\frac{HM}{\lambda}} \beta_k^2(\delta)$$

□

We are now ready to prove Proposition G.10.

Proof of Proposition G.10. Part 1. We start by writing the exact second-order Taylor expansion of the likelihood $L_\lambda^k(\mathbf{f})$, which holds for any $\mathbf{f} \succeq \mathbb{R}^{M(M+1)SAH}$

$$L_\lambda^k(\mathbf{f}) = L_\lambda^k(\hat{\mathbf{f}}_k) + r_{\mathbf{f}} L_\lambda^k(\hat{\mathbf{f}}_k)^T(\mathbf{f} - \hat{\mathbf{f}}_k) + (\mathbf{f} - \hat{\mathbf{f}}_k)^T \left(\int_{v=0}^1 (1-v) r_{\mathbf{f}}^2 L_\lambda^k(\hat{\mathbf{f}}_k + v(\mathbf{f} - \hat{\mathbf{f}}_k)) dv \right) (\mathbf{f} - \hat{\mathbf{f}}_k).$$

Since $\hat{\mathbf{f}}_k$ is the solution to the unconstrained minimization of the concave likelihood $L_\lambda^k(\mathbf{f})$, we have that $r_{\mathbf{f}} L_\lambda^k(\hat{\mathbf{f}}_k) = 0$. Recalling that $r_{\mathbf{f}}^2 L_\lambda^k(\mathbf{f}) = \mathbf{H}_k(\mathbf{f})$, we get

$$\begin{aligned} L_\lambda^k(\mathbf{f}) - L_\lambda^k(\hat{\mathbf{f}}_k) &= r_{\mathbf{f}} L_\lambda^k(\hat{\mathbf{f}}_k)^T(\mathbf{f} - \hat{\mathbf{f}}_k) + (\mathbf{f} - \hat{\mathbf{f}}_k)^T \left(\int_{v=0}^1 (1-v) r_{\mathbf{f}}^2 L_\lambda^k(\hat{\mathbf{f}}_k + v(\mathbf{f} - \hat{\mathbf{f}}_k)) dv \right) (\mathbf{f} - \hat{\mathbf{f}}_k) \\ &= (\mathbf{f} - \hat{\mathbf{f}}_k)^T \left(\int_{v=1}^1 (1-v) \mathbf{H}_k(\hat{\mathbf{f}}_k + v(\mathbf{f} - \hat{\mathbf{f}}_k)) dv \right) (\mathbf{f} - \hat{\mathbf{f}}_k) \\ &= \left\| \mathbf{f} - \hat{\mathbf{f}}_k \right\|_{\tilde{\mathbf{G}}_k(\hat{\mathbf{f}}_k, \mathbf{f})}^2 && \text{(Def. of } \tilde{\mathbf{G}}_k(\hat{\mathbf{f}}_k, \mathbf{f})) \\ &= \left\| \mathbf{f} - \hat{\mathbf{f}}_k \right\|_{\mathbf{G}_k(\hat{\mathbf{f}}_k, \mathbf{f})}^2 && (\tilde{\mathbf{G}}_k = \mathbf{G}_k) \\ &= \left\| \mathbf{g}_k(\mathbf{f}) - \mathbf{g}_k(\hat{\mathbf{f}}_k) \right\|_{\mathbf{G}_k^{-1}(\hat{\mathbf{f}}_k, \mathbf{f})}^2 && \text{(Lemma G.2)} \\ &= \left\| \mathbf{g}_k(\mathbf{f}) - \mathbf{g}_k(\hat{\mathbf{f}}_k) \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \hat{\mathbf{f}}_k)}^2. && (\mathbf{G}_k(\hat{\mathbf{f}}_k, \mathbf{f}) = \mathbf{G}_k(\mathbf{f}, \hat{\mathbf{f}}_k)) \end{aligned}$$

Rearranging, we get that for any $\mathbf{f} \succeq \mathbb{R}^{M(M+1)SAH}$,

$$L_\lambda^k(\hat{\mathbf{f}}_k) - L_\lambda^k(\mathbf{f}) \leq \left\| \mathbf{g}_k(\mathbf{f}) - \mathbf{g}_k(\hat{\mathbf{f}}_k) \right\|_{\mathbf{G}_k^{-1}(\mathbf{f}, \hat{\mathbf{f}}_k)}^2,$$

and thus, the above inequality holds for any $\mathbf{f} \succeq \mathcal{C}_k(\delta) \subseteq \mathbb{R}^{M(M+1)SAH}$.

Finally, by Lemma G.11, for any $\mathbf{f} \succeq \mathcal{C}_k(\delta)$, we have that

$$L_\lambda^k(\hat{\mathbf{f}}_k) - L_\lambda^k(\mathbf{f}) \leq \xi_k^2(\delta),$$

which implies that $\mathcal{C}_k(\delta) \subseteq E_k(\delta)$ by the definition of Equation (15). In particular, by Lemma G.5, $\mathbf{f} \succeq \mathcal{C}_k(\delta)$ with probability at least $1 - \delta$, and thus the same holds for $E_k(\delta)$.

Part 2. For this part, assume that $\mathbf{f} \succeq \mathcal{C}_k(\delta) \subseteq E_k(\delta)$ for all $k \geq 1$, an event that holds with probability at least $1 - \delta$. Also, let $\mathbf{f} \succeq E_k(\delta)$. Writing the Taylor expansion of $L_\lambda^k(\mathbf{f})$, following the same derivation as the last part, we get

$$\begin{aligned} L_\lambda^k(\mathbf{f}) &= L_\lambda^k(\mathbf{f}) + r_{\mathbf{f}} L_\lambda^k(\mathbf{f})^T(\mathbf{f} - \mathbf{f}) + (\mathbf{f} - \mathbf{f})^T \left(\int_{v=1}^1 (1-v) r_{\mathbf{f}}^2 L_\lambda^k(\mathbf{f} + v(\mathbf{f} - \mathbf{f})) dv \right) (\mathbf{f} - \mathbf{f}) \\ &= L_\lambda^k(\mathbf{f}) + r_{\mathbf{f}} L_\lambda^k(\mathbf{f})^T(\mathbf{f} - \mathbf{f}) + k_{\mathbf{f}} \mathbf{f}^T k_{\mathbf{G}_k(\mathbf{f}, \mathbf{f})}^2 \cdot \\ &= L_\lambda^k(\mathbf{f}) + r_{\mathbf{f}} L_\lambda^k(\mathbf{f})^T(\mathbf{f} - \mathbf{f}) + (2+2L)^{-1} k_{\mathbf{f}} \mathbf{f}^T k_{\mathbf{H}_k(\mathbf{f})}^2. \end{aligned} \tag{Lemma G.9}$$

Rearranging this inequality, we get,

$$\begin{aligned} k_{\mathbf{f}} \mathbf{f}^T k_{\mathbf{H}_k(\mathbf{f})}^2 &\leq (2+2L)(L_\lambda^k(\mathbf{f}) - L_\lambda^k(\hat{\mathbf{f}}_k)) + (2+2L)r_{\mathbf{f}} L_\lambda^k(\mathbf{f})^T(\mathbf{f} - \hat{\mathbf{f}}_k) \\ &= (2+2L)(L_\lambda^k(\hat{\mathbf{f}}_k) - L_\lambda^k(\mathbf{f})) + (2+2L)r_{\mathbf{f}} L_\lambda^k(\mathbf{f})^T(\mathbf{f} - \hat{\mathbf{f}}_k) && \text{(Def. of } \hat{\mathbf{f}}_k) \\ &= (2+2L)\xi_k^2(\delta) + (2+2L)r_{\mathbf{f}} L_\lambda^k(\mathbf{f})^T(\mathbf{f} - \hat{\mathbf{f}}_k) && (\mathbf{f} \succeq E_k(\delta)) \\ &= (2+2L)\xi_k^2(\delta) + (2+2L)k_{\mathbf{f}} \mathbf{f}^T k_{\mathbf{H}_k(\mathbf{f})} \left\| r_{\mathbf{f}} L_\lambda^k(\mathbf{f}) \right\|_{\mathbf{H}_k^{-1}(\mathbf{f})} && \text{(Cauchy-Schwartz)} \\ &= (2+2L)\xi_k^2(\delta) + (2+2L)\beta_k(\delta) k_{\mathbf{f}} \mathbf{f}^T k_{\mathbf{H}_k(\mathbf{f})} \end{aligned}$$

where the last inequality is since

$$\begin{aligned} \left\| r_{\mathbf{f}} L_{\lambda}^k(\mathbf{f}) \right\|_{\mathbf{H}_k^{-1}(\mathbf{f})} &= \left\| r_{\mathbf{f}} L_{\lambda}^k(\mathbf{f}) - \underbrace{r_{\mathbf{f}} L_{\lambda}^k(\hat{\mathbf{f}}_k)}_{=0} \right\|_{\mathbf{H}_k^{-1}(\mathbf{f})} \\ &= \left\| \mathbf{g}_k(\mathbf{f}) - \mathbf{g}_k(\hat{\mathbf{f}}_k) \right\|_{\mathbf{H}_k^{-1}(\mathbf{f})} \quad (\text{see Equation (11)}) \\ &= \beta_k(\delta). \quad (\mathbf{f} \geq C_k(\delta)) \end{aligned}$$

Thus, we have the inequality

$$k\mathbf{f} - \mathbf{f} \left\|_{\mathbf{H}_k(\mathbf{f})}^2 \leq (2 + 2L)\xi_k^2(\delta) + (2 + 2L)\beta_k(\delta) k\mathbf{f} - \mathbf{f} \left\|_{\mathbf{H}_k(\mathbf{f})},$$

which implies that (see, e.g. Abeille et al., 2021, Proposition 7)

$$k\mathbf{f} - \mathbf{f} \left\|_{\mathbf{H}_k(\mathbf{f})}^2 \leq \sqrt{(2 + 2L)\xi_k(\delta) + (2 + 2L)\beta_k(\delta)}.$$

To conclude the proof, notice that under the event that $\mathbf{f} \geq C_k(\delta)$ for all $k \geq 1$, we also have that $\mathbf{f} \geq E_k(\delta)$, and therefore, $E_k(\delta)$ is not empty for all $k \geq 1$. Specifically, when the set is nonempty, by the definition of the set $E_k(\delta)$ (Equation (15)), there exists a $\mathbf{f} \geq F$ for which $L_{\lambda}^k(\hat{\mathbf{f}}_k) - L_{\lambda}^k(\mathbf{f}) \leq \xi_k^2(\delta)$.

Now, by definition, it holds for the constrained maximizer $\bar{\mathbf{f}}_k \geq \arg \max_{\mathbf{f} \geq F} L_{\lambda}^k(\mathbf{f})$ that for any $\mathbf{f} \geq F$, $L_{\lambda}^k(\bar{\mathbf{f}}_k) - L_{\lambda}^k(\mathbf{f})$. Consequently, for any $\mathbf{f} \geq F$, $L_{\lambda}^k(\hat{\mathbf{f}}_k) - L_{\lambda}^k(\bar{\mathbf{f}}_k) \leq L_{\lambda}^k(\hat{\mathbf{f}}_k) - L_{\lambda}^k(\mathbf{f})$. Thus, when the set is nonempty, it must contain the constrained maximizer $\bar{\mathbf{f}}_k$. A direct conclusion of the previous inequality is that w.p. at least $1 - \delta$, for all $k \geq 1$,

$$\left\| \bar{\mathbf{f}}_k - \mathbf{f} \right\|_{\mathbf{H}_k(\mathbf{f})} \leq (2 + 2L)\beta_k(\delta) + \sqrt{2(1 + L)\xi_k(\delta)}.$$

□

G.3. Local Confidence Bound

To prove the local confidence bound, we adapt the proofs of (Tennenholtz et al., 2022, Appendix K) to the multinomial case, while also taking into account the discounting of the latent features.

Next, we prove that the inverse of the Gram matrix of each episode is well behaved – its diagonal is bounded at any visited state-action-context.

Lemma G.12 (Inverse Eigenvalues Bound). *Let $\mathbf{D}_k = \sum_{h=1}^H \mathbf{d}_h^{k^0} \mathbf{d}_h^{k^0T}$ be the Gram matrix that corresponds to the discounted visitations during episode k . If $(s, a, x) \geq \tau_h^k$ and $\mathbf{e}_{x,s,a,h} \geq \mathbb{R}^{(M+1)SAH}$ is a unit vector in the coordinate (x, s, a, h) , then $\mathbf{e}_{x,s,a,h}^T (\lambda \mathbf{I} + \mathbf{D}_k)^{-1} \mathbf{e}_{x,s,a,h} \leq \frac{1}{4H\alpha + \lambda}$.*

Proof. We closely follow the proof of (Tennenholtz et al., 2022, Lemma 7), while incorporating discount to the visitation vector. For brevity, and with some abuse of notations, we use $\mathbf{e}_n \geq \mathbb{R}^{(M+1)SAH}$ to denote the unit vector in the n -th coordinate. In the following, we assume w.l.o.g. that the t -th coordinate of the vector \mathbf{d}_h^k represents the state that was visited on the t -th time step (while unvisited states can be arbitrarily ordered). As done by (Tennenholtz et al., 2022, Lemma 7), this can be done using any permutation matrix \mathbf{P}_k such that $\mathbf{e}_{x_t^k, s_t^k, a_t^k, t} = \mathbf{P}_k \mathbf{e}_t$ for all $t \geq [H]$. Then, denoting $\bar{\mathbf{e}}_t = H_{\alpha}^{1/2} \sum_{n=1}^t \alpha^{t-n} \mathbf{e}_n = \underbrace{(\alpha^{t-1}, \alpha^{t-2}, \dots, 1, 0, \dots, 0)}_{t \text{ elements}}^T$, we can write $\mathbf{d}_t^k = H_{\alpha}^{1/2} \sum_{n=1}^t \alpha^{t-n} \mathbf{e}_{x_n^k, s_n^k, a_n^k, n} =$

$$H_{\alpha}^{1/2} \sum_{n=1}^t \alpha^{t-n} \mathbf{P}_k \mathbf{e}_n = \mathbf{P}_k \bar{\mathbf{e}}_t.$$

Now, recalling that permutation matrices are orthogonal ($\mathbf{P}_k^{-1} = \mathbf{P}_k^T$) we can write

$$\begin{aligned}
 \mathbf{e}_{x,s,a,h}^T (\lambda \mathbf{I} + \mathbf{D}_k)^{-1} \mathbf{e}_{x,s,a,h} &= \mathbf{e}_{x,s,a,h}^T \left(\lambda \mathbf{I} + \sum_{t=1}^H \mathbf{d}_t^{k^0} \mathbf{d}_t^{k^0 T} \right)^{-1} \mathbf{e}_{x,s,a,h} \\
 &= \mathbf{e}_{x,s,a,h}^T \left(\lambda \mathbf{I} + \sum_{t=1}^H \mathbf{P}_k \bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T \mathbf{P}_k^T \right)^{-1} \mathbf{e}_{x,s,a,h} \\
 &= \mathbf{e}_{x,s,a,h}^T \left(\mathbf{P}_k \left(\lambda \mathbf{I} + \sum_{t=1}^H \bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T \right) \mathbf{P}_k^T \right)^{-1} \mathbf{e}_{x,s,a,h} \\
 &= \mathbf{e}_{x,s,a,h}^T \mathbf{P}_k \left(\lambda \mathbf{I} + \sum_{t=1}^H \bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T \right)^{-1} \mathbf{P}_k^T \mathbf{e}_{x,s,a,h} \\
 &= \mathbf{e}_h^T \left(\lambda \mathbf{I} + \sum_{t=1}^H \bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T \right)^{-1} \mathbf{e}_h .
 \end{aligned}$$

Next, notice that $\lambda \mathbf{I} + \sum_{t=1}^H \bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T$ is a block-diagonal matrix, whose first block is of size $H \times H$ (and the rest of the matrix is fully diagonal). We denote this first block of $\sum_{t=1}^H \bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T$ by \mathbf{C} . For block-diagonal matrices, each block can be inverted independently of the other blocks, and for any coordinate $h \in [H]$, if $\mathbf{u}_h \in \mathbb{R}^H$ is the unit vector at coordinate h , we thus have

$$\begin{aligned}
 \mathbf{e}_{x,s,a,h}^T (\lambda \mathbf{I} + \mathbf{D}_k)^{-1} \mathbf{e}_{x,s,a,h} &= \mathbf{e}_h^T \left(\lambda \mathbf{I} + \sum_{t=1}^H \bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T \right)^{-1} \mathbf{e}_h \\
 &= \mathbf{u}_h^T (\lambda \mathbf{I} + \mathbf{C})^{-1} \mathbf{u}_h \\
 &= \frac{k \mathbf{u}_h k_2^2}{\lambda_{\min}(\lambda \mathbf{I} + \mathbf{C})} \\
 &= \frac{1}{\lambda + \lambda_{\min}(\mathbf{C})}
 \end{aligned} \tag{16}$$

In the rest of the proof, we focus on bounding $\lambda_{\min}(\mathbf{C})$. First, observe that for any $t \in [H]$, we have

$$(\bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T)(i, j) = \begin{cases} H \alpha^{-1} \alpha^{2t} & i = j = t \\ 0 & \text{else} \end{cases}$$

and thus

$$\mathbf{C}(i, j) = \sum_{t=1}^H (\bar{\mathbf{e}}_t \bar{\mathbf{e}}_t^T)(i, j) = H \alpha^{-1} \sum_{t=\max\{\bar{r}_i, \bar{r}_j\}}^H \alpha^{2t} & i = j \\
 & 0 & \text{else}$$

In particular, notice that for any $i < j$ (above diagonal), we have $\mathbf{C}(i, j) = \alpha \mathbf{C}(i+1, j)$, while for $i = j$ (below and on

diagonal), we have $C(i, j) = \alpha C(i + 1, j) + \alpha^{i-j}$. Using this structure, we can calculate its inverse using diagonalization:

$$\begin{aligned}
 & \left(\begin{array}{cccc|cccc}
 \sum_{t=1}^H \alpha^{2t-2} & \sum_{t=2}^H \alpha^{2t-3} & \sum_{t=3}^H \alpha^{2t-4} & \dots & \alpha^{H-1} & 1 & 0 & 0 & \dots & 0 \\
 \sum_{t=2}^H \alpha^{2t-3} & \sum_{t=2}^H \alpha^{2t-4} & \sum_{t=3}^H \alpha^{2t-5} & \dots & \alpha^{H-2} & 0 & 1 & 0 & \dots & 0 \\
 \sum_{t=3}^H \alpha^{2t-4} & \sum_{t=3}^H \alpha^{2t-5} & \sum_{t=3}^H \alpha^{2t-6} & \dots & \alpha^{H-3} & 0 & 0 & 1 & \dots & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 \alpha^{H-1} & \alpha^{H-2} & \alpha^{H-3} & \dots & 1 & 0 & 0 & 0 & \dots & 1
 \end{array} \right) \\
 &= \left(\begin{array}{cccc|cccc}
 1 & 0 & 0 & \dots & 0 & 1 & \alpha & 0 & \dots & 0 \\
 \alpha & 1 & 0 & \dots & 0 & 0 & 1 & \alpha & \dots & 0 \\
 \alpha^2 & \alpha & 1 & \dots & 0 & 0 & 0 & 1 & \dots & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 \alpha^{H-1} & \alpha^{H-2} & \alpha^{H-3} & \dots & 1 & 0 & 0 & 0 & \dots & 1
 \end{array} \right) \\
 &= \left(\begin{array}{cccc|cccc}
 1 & 0 & 0 & \dots & 0 & 1 & \alpha & 0 & \dots & 0 \\
 0 & 1 & 0 & \dots & 0 & \alpha & 1 + \alpha^2 & \alpha & \dots & 0 \\
 0 & 0 & 1 & \dots & 0 & 0 & \alpha & 1 + \alpha^2 & \dots & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & \dots & 1 + \alpha^2
 \end{array} \right)
 \end{aligned}$$

In the first relation, we subtracted α -times the $i + 1$ rows from the i rows, while in the second one, we subtracted α -times the $i - 1$ rows from the i rows. Thus, the inverse can be explicitly written as:

$$C_{i,j}^{-1} = H_\alpha \begin{cases} 1 & i = j = 1 \\ 1 + \alpha^2 & i = j > 1 \\ \alpha & i = j - 1 \text{ or } i = j + 1 \\ 0 & \text{o.w.} \end{cases}$$

Notice that the absolute values of all rows is smaller than $H_\alpha(1 + \alpha)^2 = 4H_\alpha$. Then (e.g., by Gershgorin circle theorem), $\lambda_{\max}(C^{-1}) \leq 4H_\alpha$, and since B is PSD, $\lambda_{\min}(C) \geq \frac{1}{4H_\alpha}$. The proof is concluded by substituting this result back into Equation (16). \square

We are now ready to prove the local concentration results for the latent features of a DCM DP:

Lemma G.13 (Local Estimation Confidence Bound). *Let $\hat{\mathbf{f}}^k \in \arg \max_{\mathbf{f} \in \mathcal{F}} L_\lambda^k(\mathbf{f})$ be the maximum likelihood estimate of the features. Then, for any $\delta > 0$, with probability of at least $1 - \delta$, for all $k \in [K]$, $h \in [H]$, $i \in [M]$ and $s, a, x \in \mathcal{S} \times \mathcal{A} \times \mathcal{X}$, it holds that*

$$\left| \hat{f}_{i,h}^k(s, a, x) - f_{i,h}(s, a, x) \right| \leq \frac{2\gamma_k(\delta) \sqrt{\kappa H_\alpha}}{\sqrt{n_h^k(s, a, x) + 4\lambda H_\alpha}},$$

where $\gamma_k(\delta)$ is defined in Proposition G.10.

Proof. The proof follows Lemma 6 of (Tennenholtz et al., 2022).

For any $k \in [K]$, $h \in [H]$, $i \in [M]$ and $s, a, x \in \mathcal{S} \times \mathcal{A} \times \mathcal{X}$, let $\mathbf{e}_{x,s,a,h} \in \mathbb{R}^{(M+1)SAH}$ be a unit vector in the (x, s, a, h) coordinate and denote $\mathbf{f}_i \in \mathbb{R}^{(M+1)SAH}$ the latent features that correspond to a next latent state i . We start by bounding

$$\begin{aligned}
 \left| \hat{f}_{i,h}^k(s, a, x) - f_{i,h}(s, a, x) \right| &= \left| \langle \mathbf{e}_{x,s,a,h}, \hat{\mathbf{f}}_i - \mathbf{f}_i \rangle \right| \\
 &\leq k \mathbf{e}_{x,s,a,h}^T \mathbf{V}_k^{-1} \left\| \hat{\mathbf{f}}_i - \mathbf{f}_i \right\|_{\mathbf{V}_k} \quad (\text{Cauchy-Schwartz}) \\
 &\leq k \mathbf{e}_{x,s,a,h}^T \mathbf{V}_k^{-1} \sqrt{\sum_{i^0=1}^M \left\| \hat{\mathbf{f}}_{i^0} - \mathbf{f}_{i^0} \right\|_{\mathbf{V}_k}^2} \\
 &= k \mathbf{e}_{x,s,a,h}^T \mathbf{V}_k^{-1} \left\| \hat{\mathbf{f}} - \mathbf{f} \right\|_{\mathbf{I}_M \otimes \mathbf{V}_k}, \quad (17)
 \end{aligned}$$

We now turn our focus to bound $ke_{x,s,a,h}k_{\mathbf{V}_k}^{-1}$. Using the notation \mathbf{D}_k , as defined in Lemma G.12, we have

$$\begin{aligned} \mathbf{V}_k^{-1} &= \left(\lambda \mathbf{I} + \sum_{k^0=1}^k \mathbf{D}_{k^0} \mathbf{1}\{(s, a, x) \succeq \tau_h^{k^0}\} \right)^{-1} \\ &= \left(\sum_{k^0=1}^k \left(\frac{\lambda}{n_h(s, a, x)} \mathbf{I} + \mathbf{D}_{k^0} \right) \mathbf{1}\{(s, a, x) \succeq \tau_h^{k^0}\} \right)^{-1} \\ &= \frac{1}{(n_h(s, a, x))^2} \sum_{k^0=1}^k \left(\frac{\lambda}{n_h(s, a, x)} \mathbf{I} + \mathbf{D}_{k^0} \right)^{-1} \mathbf{1}\{(s, a, x) \succeq \tau_h^{k^0}\}, \end{aligned}$$

where $n_h^k(s, a, x) = \sum_{k^0=1}^k \mathbf{1}\{(s, a, x) \succeq \tau_h^{k^0}\}$ and the third transition is due to HM-AM inequality for positive matrices (Bhagwat and Subramanian, 1978). Next, we combine this result with Lemma G.12 and get

$$\begin{aligned} ke_{x,s,a,h}k_{\mathbf{V}_k}^2 &= e_{x,s,a,h}^T \mathbf{V}_k^{-1} e_{x,s,a,h} \\ &= \frac{1}{(n_h^k(s, a, x))^2} \sum_{k^0=1}^k e_{x,s,a,h}^T \left(\frac{\lambda}{n_h(s, a, x)} \mathbf{I} + \mathbf{D}_{k^0} \right)^{-1} e_{x,s,a,h} \mathbf{1}\{(s, a, x) \succeq \tau_h^{k^0}\} \\ &= \frac{1}{(n_h(s, a, x))^2} \sum_{k^0=1}^k \frac{1}{\frac{1}{4H_\alpha} + \frac{\lambda}{n_h(s, a, x)}} \mathbf{1}\{(s, a, x) \succeq \tau_h^{k^0}\} \quad (\text{Lemma G.12}) \\ &= \frac{n_h(s, a, x)}{(n_h(s, a, x))^2} \frac{1}{\frac{1}{4H_\alpha} + \frac{\lambda}{n_h(s, a, x)}} \\ &= \frac{1}{\frac{n_h(s, a, x)}{4H_\alpha} + \lambda} \\ &= \frac{4H_\alpha}{n_h(s, a, x) + 4\lambda H_\alpha}. \end{aligned}$$

By plugging into Equation (17), we obtain that for any k and any h, s, a

$$\begin{aligned} \left| \hat{f}_{i,h}^k(s, a, x) - f_{i,h}(s, a, x) \right| &\leq ke_{x,s,a,h}k_{\mathbf{V}_k}^{-1} \left\| \hat{\mathbf{f}} - \mathbf{f} \right\|_{\mathbf{I}_M \rho_{\mathbf{V}_k}^{\mathbf{V}_k}} \\ &= \left\| \hat{\mathbf{f}} - \mathbf{f} \right\|_{\mathbf{I}_M \mathbf{V}_k} \frac{2 \rho_{\mathbf{H}_\alpha}^{\mathbf{V}_k}}{\sqrt{n_h(s, a, x) + 4\lambda H_\alpha}} \\ &= \left\| \hat{\mathbf{f}} - \mathbf{f} \right\|_{\mathbf{H}_k(\mathbf{f})} \frac{2 \rho_{\kappa H_\alpha}^{\mathbf{V}_k}}{\sqrt{n_h(s, a, x) + 4\lambda H_\alpha}}. \end{aligned}$$

Finally by Proposition G.10, with probability $1 - \delta$, for all $k \geq 1$, it holds that $\left\| \hat{\mathbf{f}} - \mathbf{f} \right\|_{\mathbf{H}_k(\mathbf{f})} \leq \gamma_k(\delta)$, and substituting this bound concludes the proof. \square