

Sample-efficient Learning of Infinite-horizon Average-reward MDPs with General Function Approximation

Jianliang He^{*} Han Zhong[†] Zhuoran Yang[‡]

Abstract

We study infinite-horizon average-reward Markov decision processes (AMDPs) in the context of general function approximation. Specifically, we propose a novel algorithmic framework named Local-fitted Optimization with OPTimism (LOOP), which incorporates both model-based and value-based incarnations. In particular, LOOP features a novel construction of confidence sets and a low-switching policy updating scheme, which are tailored to the average-reward and function approximation setting. Moreover, for AMDPs, we propose a novel complexity measure — average-reward generalized eluder coefficient (AGEC) — which captures the challenge of exploration in AMDPs with general function approximation. Such a complexity measure encompasses almost all previously known tractable AMDP models, such as linear AMDPs and linear mixture AMDPs, and also includes newly identified cases such as kernel AMDPs and AMDPs with Bellman eluder dimensions. Using AGEC, we prove that LOOP achieves a sub-linear $\tilde{O}(\text{poly}(d, \text{sp}(V^*))\sqrt{T\beta})$ regret, where d and β correspond to AGEC and log-covering number of the hypothesis class respectively, $\text{sp}(V^*)$ is the span of the optimal state bias function, T denotes the number of steps, and $\tilde{O}(\cdot)$ omits logarithmic factors. When specialized to concrete AMDP models, our regret bounds are comparable to those established by the existing algorithms designed specifically for these special cases. To the best of our knowledge, this paper presents the first comprehensive theoretical framework capable of handling nearly all AMDPs.

Contents

1	Introduction	3
2	Preliminaries	4
3	General Function Approximation	6
3.1	Average-Reward Generalized Eluder Coefficients	7
3.2	Relation with Tractable Complexity Metric	9

^{*}Fudan University. Email: hejl20@fudan.edu.cn.

[†]Peking University. Email: hanzhong@stu.pku.edu.cn.

[‡]Yale University. Email: zhuoran.yang@yale.edu.

4	Local-fitted Optimization with Optimism	11
5	Proof Overview of Regret Analysis	12
6	Conclusion	14
A	Backgrounds and Technical Novelties	19
B	Discussions	21
C	Alternative Choices of Discrepancy Function	22
D	Concrete Examples	23
D.1	Linear Function Approximation and Variants	23
D.2	Kernel Function Approximation	25
D.3	Linear Mixture AMDP	26
E	Proof of Main Results for LOOP	27
E.1	Proof of Theorem 4.1	27
E.2	Proof of Lemma E.1	28
E.3	Proof of Lemma E.2	30
E.4	Proof of Lemma E.3	31
F	Proof of Results about Complexity Measures	34
F.1	Proof of Lemma 3.1	34
F.2	Proof of Lemma 3.2	35
G	Proof of Results for Concrete Examples	35
G.1	Proof of Proposition D.1	35
G.2	Proof of Proposition D.2	37
G.3	Proof of Proposition D.3	37
G.4	Proof of Proposition D.4	38
G.5	Discussion about Performance on Concrete Examples	40
H	Technical Lemmas	42
H.1	Proof of Technical Lemmas	44
H.1.1	Proof of Lemma H.2	44
H.1.2	Proof of Lemma H.3	44
H.1.3	Proof of Lemma H.4	46
I	Supplementary Discussions	46
I.1	Proof Sketch of MLE-based Results	46
I.2	Extended Value Iteration (EVI) for Model-Based Hypotheses	48

1 Introduction

Reinforcement learning (RL) (Sutton and Barto, 2018) is a powerful tool for addressing intricate sequential decision-making problems. In this context, Markov decision processes (MDPs) (Puterman, 2014; Sutton and Barto, 2018) frequently serve as a fundamental model for modeling such decision-making scenarios. Motivated by different feedback structures in real applications, MDPs consist of three subclasses — finite-horizon MDPs, infinite-horizon discounted MDPs, and infinite-horizon average-reward MDPs. Each of these MDP variants is of paramount significance and operates in a parallel fashion, with none being amenable to complete reduction into another. Of these three MDP subclasses, finite-horizon MDPs have received significant research efforts in understanding their exploration challenge, especially in the presence of large state spaces which necessitates function approximation tools. Existing works on finite-horizon MDPs have proposed numerous structural conditions on the MDP model that empower sample-efficient learning. These structural conditions include but are not limited to linear function approximation (Jin et al., 2020), Bellman rank (Jiang et al., 2017), eluder dimension (Wang et al., 2020), Bellman eluder dimension (Jin et al., 2021), bilinear class (Du et al., 2021), decision estimation coefficient (Foster et al., 2021), and generalized eluder coefficient (Zhong et al., 2022). Moreover, these works have designed various model-based and value-based algorithms to address finite-horizon MDPs governed by these structural conditions.

In contrast to the rich literature devoted to finite-horizon MDPs, the study of sample-efficient exploration in infinite-horizon MDPs has hitherto been relatively limited. Importantly, it remains elusive how to design in a principled fashion a sample-efficient RL algorithm in the online setting with general function approximation. To this end, we focus on infinite-horizon average-reward MDPs (AMDPs), which offer a suitable framework for addressing real-world decision-making scenarios that prioritize long-term returns, such as product delivery (Proper and Tadepalli, 2006). Our work endeavors to provide a unified theoretical foundation for understanding infinite-horizon average-reward MDPs from the perspective of general function approximation, akin to the comprehensive investigations conducted in the domain of finite-horizon MDPs. To pursue this overarching objective, we have delineated two subsidiary goals that form the crux of our research endeavor.

- **Development of a Novel Structural Condition/Complexity Measure.** Existing works are restricted to tabular AMDPs (Bartlett and Tewari, 2012; Jaksch et al., 2010) and AMDPs with linear function approximation (Wu et al., 2022; Wei et al., 2021), with Chen et al. (2022a) as the only exception (to our best knowledge). While Chen et al. (2022b) does extend the eluder dimension for finite-horizon MDPs (Ayoub et al., 2020) to the infinite-horizon average-reward context, their complexity measure seems to be only slightly more general than the linear mixture AMDPs (Wu et al., 2022) and falls short in capturing other fundamental models such as linear AMDPs (Wei et al., 2021). Hence, our first subgoal is proposing a new structural condition. This condition is envisioned to be sufficiently versatile to encompass all known tractable AMDPs, while also potentially introducing innovative and tractable models into the framework.
- **Algorithmic Framework for Addressing Identified Structural Condition.** The second subgoal is anchored in the development of sample-efficient algorithms for AMDPs characterized by the structural condition proposed in our work. Our aspiration is to devise an algorithmic framework that can be flexibly implemented in both model-based and value-based paradigms,

depending on the nature of the problem at hand. This adaptability guarantees that our algorithms possess the ability to effectively address a wide range of AMDPs.

Our work attains these two pivotal subgoals through the introduction of (i) a novel complexity measure — Average-reward Generalized Eluder Coefficient (AGEC), and (ii) a corresponding algorithmic framework dubbed as Local-fitted Optimization with OPTimism (LOOP). Our primary contributions and novelties are summarized below:

- **AGEC Complexity Measure.** Our complexity measure AGEC extends the generalized eluder coefficient (GEC) in [Zhong et al. \(2022\)](#) to the infinite-horizon average-reward setting. However, it incorporates significant modifications. AGEC not only establishes a connection between the Bellman error and the training error, akin to GEC but also imposes certain constraints on transferability (see Definition 3 for details). This modification proves instrumental in attaining sample efficiency in the realm of AMDPs (see Section 5 for detailed discussion). We demonstrate that AGEC not only encompasses all previously recognized tractable AMDPs, including tabular AMDPs ([Bartlett and Tewari, 2012](#); [Jaksch et al., 2010](#)), linear AMDPs ([Wei et al., 2021](#)), linear mixture MDPs ([Wu et al., 2022](#)), AMDPs with low eluder dimension ([Chen et al., 2022a](#)), but also captures some new identified models like linear Q^*/V^* AMDPs (see Definition 13), kernel AMDPs (see Proposition D.3), and AMDPs with low Bellman eluder dimension (see Definition 8).
- **LOOP Algorithmic Framework.** Our algorithm LOOP is based on the optimism principle and features a novel construction of confidence sets along with a low-switching updating scheme. Remarkably, LOOP offers the flexibility to be implemented either in the model-based or value-based paradigm, depending on the problem type.
- **Unified Theoretical Results.** From the theoretical side, we prove that LOOP enjoys the regret of $\tilde{\mathcal{O}}(\text{poly}(d, \text{sp}(V^*))\sqrt{T\beta})$, where d and β correspond to the AGEC and the log-covering number of the hypothesis class respectively, $\text{sp}(V^*)$ denotes the span of the optimal state bias function, T is the number of steps, and $\tilde{\mathcal{O}}$ hides logarithmic factors. This result shows that LOOP is capable of solving all AMDPs with low AGEC.

In summary, we provide a unified theoretical understanding of infinite-horizon AMDPs with general function approximation. Further elaboration on our contributions and technical novelties are provided in Appendix A. Due to space limits, we only provide a comparison between our results and mostly related works on AMDPs in Table 1. More related works are deferred to Appendix A.

2 Preliminaries

Notations. For any integer $n \in \mathbb{N}^+$, we take the convention to use $[n] = \{1, \dots, n\}$. Consider two non-negative sequences $\{a_n\}_{n \geq 0}$ and $\{b_n\}_{n \geq 0}$, if $\limsup a_n/b_n < \infty$, then we write it as $a_n = \mathcal{O}(b_n)$. Else if $\limsup a_n/b_n = 0$, then we write it as $a_n = o(b_n)$. And we use $\tilde{\mathcal{O}}$ to omit the logarithmic terms. Denote $\Delta(\mathcal{X})$ be the probability simplex over the set \mathcal{X} . Denote by $\sup_x |v(x)|$ the supremum norm of a given function. $x \wedge y$ stands for $\min\{x, y\}$ and $x \vee y$ stands for $\max\{x, y\}$. Given any continuum \mathcal{S} , let $|\mathcal{S}|$ be the cardinality. Given two distributions $P, Q \in \Delta(\mathcal{X})$, the TV distance of the two distributions is defined as $\text{TV}(P, Q) = \frac{1}{2} \mathbb{E}_{x \sim P}[|dQ(x)/dP(x) - 1|]$.

Algorithm	Assumption	Type	Tabular	Linear Mixture	Linear	Eluder	ABE	Kernel	AGEC
LOOP (Ours)	Bellman optimality (finite span)	Model-based & Value-based	✓	✓	✓	✓	✓	✓	✓
SIM-TO-REAL (Chen et al., 2022a)	Communicating AMDP (finite diameter)	Model-based	✓	✓	✗	✓	✗	✗	✗
UCRL2-VTR (Wu et al., 2022)	Communicating AMDP (finite diameter)	Model-based	✓	✓	✗	✗	✗	✗	✗
FOPO (Wei et al., 2021)	Bellman optimality (finite span)	Value-based	✓	✗	✓	✗	✗	✗	✗
UCRL2 (Jaksch et al., 2010)	Communicating AMDP (finite diameter)	Model-based	✓	✗	✗	✗	✗	✗	✗

Table 1: A comparison with the most related algorithms on AMDPs. We remark that our assumption is weaker since the communicating MDP satisfies the Bellman optimality and the diameter is bound by the span. Besides, average-reward Bellman eluder dimension (ABE), kernel AMDPs, and AGEK are new complexity measures proposed by our work. In particular, AGEK serves as a unifying complexity measure capable of encompassing all other established complexity measures.

An infinite-horizon average-reward Markov Dependent Process (AMDPs) is characterized by $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, \mathbb{P})$, where \mathcal{S} is a Borel state space with a possibly infinite number of elements, \mathcal{A} is a finite set of actions, $r : \mathcal{S} \times \mathcal{A} \mapsto [-1, 1]$ is an unknown reward function¹ and $\mathbb{P}(\cdot|s, a)$ is the unknown transition kernel. The learning protocol for infinite-horizon average-reward RL is as follows: the agent interacts with \mathcal{M} over a fixed number of T steps, starting from a pre-determined initial state $s_1 \in \mathcal{S}$. At each step $t \in [T]$, the agent observe a state $s_t \in \mathcal{S}$ and takes an action $a_t \in \mathcal{A}$, receives a reward $r(s_t, a_t)$ and transits to the next state s_{t+1} drawn from $\mathbb{P}(\cdot|s_t, a_t)$.

Denote $\Delta(\mathcal{A})$ be the probability simplex over the action space \mathcal{A} . Specifically, the stationary policy π is a mapping $\pi : \mathcal{S} \mapsto \Delta(\mathcal{A})$ with $\pi(a|s)$ specifying the probability of taking action a at state s . Given a stationary policy π , the long-term average reward starting is defined as

$$J^\pi(s) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) | s_1 = s \right], \quad \forall s \in [\mathcal{S}],$$

where the expectation is taken with respect to the policy π and transition \mathbb{P} . In the infinite-horizon average-reward RL, existing works mostly rely on additional assumptions to achieve sample efficiency. The necessity arises from the absence of a natural counterpart to the celebrated Bellman optimality equation in the average-reward RL that is self-evident and crucial within episodic and discounted settings (Puterman, 2014). To this end, we consider a broad subclass where a modified version of the Bellman optimality equation holds (Hernández-Lerma, 2012).

Assumption 1 (Bellman optimality equation). There exists bounded measurable function $Q^* :$

¹Throughout this paper, we consider the deterministic reward for notational simplicity and all results are readily generalized to the stochastic setting. Also, we assume reward lies in $[-1, 1]$ without loss of generality.

$\mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$, $V^* : \mathcal{S} \mapsto \mathbb{R}$ and unique constant $J^* \in [-1, 1]$ such that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, it holds

$$J^* + Q^*(s, a) = r(s, a) + \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a)}[V^*(s')], \quad V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a). \quad (2.1)$$

The Bellman optimality equation, adapted for average-reward RL, posits that for any initial states $s_1 \in \mathcal{S}$, the optimal reward is independent such that $J^*(s_1) = J^*$ under a deterministic optimal policy with $\pi^*(\cdot) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(\cdot, a)$. The justification is presented in [Wei et al. \(2021\)](#).

Note that functions $V^*(s)$ and $Q^*(s, a)$ reveal the relative advantage of starting from state s and state-action pair (s, a) under the optimal policy, and are respectively called the optimal state and state-action bias function ([Wei et al., 2021](#)). Denote $\operatorname{sp}(V) = \sup_{s, s' \in \mathcal{S}} |V(s) - V(s')|$ as the span of any bounded measurable function. Note that for any solution pair (V^*, Q^*) satisfying the Bellman optimality equation in (2.1), the shifted pair $(V^* - c, Q^* - c)$ for any constant c is still a solution. Thus, without loss of generality, we can focus on the unique centralized solution such that $\|V^*\|_\infty \leq \frac{1}{2}\operatorname{sp}(V^*)$. Following the tradition in the average-reward RL ([Wei et al., 2020, 2021](#); [Wang et al., 2022](#); [Zhang and Xie, 2023](#)), the span $\operatorname{sp}(V^*)$ is assumed to be known.

As aforementioned in the paper, distinct assumptions have been employed in average-reward RL research to ensure the explorability of the problem, which includes ergodic AMDPs ([Wei et al., 2020](#); [Hao et al., 2021](#); [Zhang and Xie, 2023](#)), communicating AMDPs ([Chen et al., 2022a](#); [Wang et al., 2022](#); [Wu et al., 2022](#)) and the Bellman optimality equation ([Wei et al., 2021](#)). Among these widely adopted assumptions, we remark that the Bellman optimality equation is the least stringent one. Note that the ergodic MDP suggests the existence of bias functions for each $\pi \in \Pi$, while the latter two only require the existence of bias functions for the optimal policy. As for weak communicating assumption, a weaker form of communicating MDP ([Wang et al., 2022](#)), it directly implies the existence of the Bellman optimality equation and thus is stronger ([Hernández-Lerma, 2012](#)). Given the Bellman optimality assumption in (2.1), we introduce the average-reward Bellman operator below:

$$(\mathcal{T}_J F)(s, a) := r(s, a) + \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a)} \left[\max_{a' \in \mathcal{A}} F(s', a') \right] - J, \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad (2.2)$$

for any bounded function $F : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ and constant $J \in [-1, 1]$. Then, the Bellman optimality equation in (2.1) can be written as $\mathcal{T}_{J^*} Q^* = Q^*$. Moreover, we define the Bellman error:

$$\mathcal{E}(F, J)(s, a) := F(s, a) - (\mathcal{T}_J F)(s, a), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}. \quad (2.3)$$

Learning Objective Under the framework of online learning for AMDPs, the agent aims to learn the optimal policy by interacting with the environment over potentially infinite steps. The (empirical) regret measures the cumulative difference between the optimal average-reward and the reward achieved after interacting for T steps, formally defined as $\operatorname{Reg}(T) = \sum_{t=1}^T (J^* - r(s_t, a_t))$.

3 General Function Approximation

To capture both model-free and model-based problems with function approximation, we consider a general hypotheses class \mathcal{H} which contains a class of functions. We consider two kinds of hypothesis

classes, targeting at *value-based* problems and *model-based* problems respectively.

Definition 1 (Value-based hypothesis). We say \mathcal{H} is a value-based hypotheses class if all hypothesis $f \in \mathcal{H}$ is defined over state-action bias function Q and average-reward J such that $f = (Q_f, J_f) \in \mathcal{H}$. Let $V_f(\cdot) = \max_{a \in \mathcal{A}} Q_f(\cdot, a)$ and $\pi_f(\cdot) = \operatorname{argmax}_{a \in \mathcal{A}} Q_f(\cdot, a)$ be the greedy bias function and policy induced from hypothesis $f \in \mathcal{H}$. Denote f^* be the optimal hypothesis under true model \mathcal{M} .

Definition 2 (Model-based hypothesis). We say \mathcal{H} is a model-based hypotheses class if all hypothesis $f \in \mathcal{H}$ is defined over the transition kernel \mathbb{P} and reward function r such that $f = (\mathbb{P}_f, r_f) \in \mathcal{H}$. Let Q_f , V_f , J_f , and π_f respectively be the optimal bias functions, average-reward and policy induced from hypothesis $f \in \mathcal{H}$, which satisfies the Bellman optimality equation such that

$$Q_f(s, a) + J_f = (r_f + \mathbb{P}_f V_f)(s, a), \quad \mathbb{P}_{f'} V_f(s, a) := \mathbb{E}_{s' \sim \mathbb{P}_{f'}(\cdot | s, a)} [V_f(s')],$$

for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$. Denote f^* as the hypothesis concerning the true model \mathcal{M} .

The definition of hypotheses class \mathcal{H} over the value-based (see Definition 1) and the model-based (see Definition 2) problems in AMDP is different from the episodic setting (Du et al., 2021; Zhong et al., 2022). The most significant difference is that the Bellman equation has a different form. As a result, in the value-based scenario, instead of using a single state-action value function Q_f in episodic setting, the paired hypothesis (Q_f, J_f) is introduced to fully capture the average-reward structure. Besides, we retain the definition of hypothesis over model-based problems, augmenting it with an additional average-reward term J_f induced from (\mathbb{P}_f, r_f) . Since we do not impose any specific structural form to the hypothesis class, we stay in the realm of *general function approximation*.

As function approximation is challenging without further assumptions (Krishnamurthy et al., 2016), we introduce the realizability assumption, which is widely adopted (Jin et al., 2021).

Assumption 2 (Realizability). We assume that $f^* \in \mathcal{H}$.

Moreover, we establish the fundamental distribution families over the state-action pair upon which the metric is built. Considering the learning goal defined over the empirical regret, throughout the paper we focus on the point-wise distribution family $\mathcal{D}_\Delta = \{\delta_{s,a}(\cdot) | (s, a) \in \mathcal{S} \times \mathcal{A}\}$, which includes collections of Dirac probability measure over $\mathcal{S} \times \mathcal{A}$. Discussions are deferred to Appendix B.

3.1 Average-Reward Generalized Eluder Coefficients

In this subsection, we are going to introduce a novel metric — average-reward generalized eluder coefficients (AGEC), to capture the complexity of hypothesis class \mathcal{H} for AMDP. Extended from the generalized Eluder coefficients (GEC; Zhong et al., 2022) for finite-horizon MDPs, AGECE is a variant to fit the infinite-horizon learning with average reward, and imposes an additional structural constraint — transferability.

Definition 3 (AGEC). Given hypothesis class \mathcal{H} , discrepancy function set $\{l_f\}_{f \in \mathcal{H}}$ and constant $\epsilon > 0$, the average-reward generalized eluder coefficients $\text{AGEC}(\mathcal{H}, \{l_f\}, \epsilon)$ is defined as the smallest coefficients κ_G and d_G , such that following two conditions hold with absolute constants $C_1, C_2 > 0$:

(i) (Bellman dominance) There exists constant $d_G > 0$ such that

$$\underbrace{\sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t)}_{\text{Bellman error}} \leq \underbrace{\left[d_G \cdot \sum_{t=1}^T \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \right]^{1/2}}_{\text{In-sample training error}} + \underbrace{C_1 \cdot \text{sp}(V^*) \min\{d_G, T\} + T\epsilon}_{\text{Burn-in cost}}.$$

(ii) (Transferability) There exists constant $\kappa_G > 0$ such that for hypotheses $f_1, \dots, f_T \in \mathcal{H}$, if $\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \leq \beta$ holds for all $t \in [T]$, then we have

$$\underbrace{\sum_{t=1}^T \|\mathbb{E}_{\zeta_t}[l_{f_t}(f_t, f_t, \zeta_t)]\|_2^2}_{\text{Out-sample training error}} \leq \kappa_G \cdot \beta \log T + \underbrace{C_2 \cdot \text{sp}(V^*)^2 \min\{\kappa_G, T\} + 2T\epsilon^2}_{\text{Burn-in cost}}.$$

In the definition above, ζ_i is a subset of trajectory with varying meaning concerning the specific choice of discrepancy function, and the expectation is taken over it; C_1, C_2 are absolute constants related to $\text{span } \text{sp}(V^*)$. To simplify the notation, we denote $\mathcal{E}(f_t)(s, a) := \mathcal{E}(Q_{f_t}, J_{f_t})(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Besides, the Burn-in cost is taken at the worst case and it varies across different settings but usually non-dominating. The intuition behind the metric is that, on average, if hypotheses have small in-sample training error on the well-explored dataset, then the prediction error on a different trajectory is expected to maintain a consistently low level (Zhong et al., 2022). In specific, the dominance coefficient d_G encapsulates the challenge inherent in assessing the performance of prediction, specifically the Bellman error, given the consistently controlled in-sample training error within the designated function class \mathcal{H} . Moreover, due to the unique challenge of infinite-horizon average-reward setting, we introduce the transferability coefficient κ_G to quantify the transferability from the in-sample training error to the out-of-sample ones. Despite this additional structural condition, we can verify that nearly all tractable AMDPs admit a low AGECE value (see Section 3.2). Moreover, in Section 5, we will demonstrate the importance of such additional structural conditions for achieving sample efficiency in AMDPs from the theoretical perspective.

Moreover, to facilitate further theoretical analysis, we make further assumptions on the discrepancy function and hypothesis class as Chen et al. (2022b); Zhong et al. (2022).

Assumption 3 (Boundedness). Given any $f \in \mathcal{H}$, it holds that $\|l_f\|_\infty \leq C \cdot \text{sp}(V^*)$ with $C > 0$.

The boundedness assumption is reasonable and uniformly satisfied, as in most cases, it takes the Bellman discrepancy, defined as: for all $\zeta_t = \{s_t, a_t, r_t, s_{t+1}\} \in \mathcal{S} \times \mathcal{A} \times \mathbb{R} \times \mathcal{S}$, we have

$$l_{f'}(f, g, \zeta_t) = Q_g(s_t, a_t) - r(s_t, a_t) - V_f(s_{t+1}) + J_g, \quad (3.1)$$

or other natural derivatives, so that the discrepancy is generally upper bounded by $\mathcal{O}(\text{sp}(V^*))$.

Assumption 4 (Generalized completeness). Let \mathcal{G} be an auxiliary function class and there exists a functional operator $\mathcal{P} : \mathcal{H} \mapsto \mathcal{G}$, we say that \mathcal{H} satisfies generalized completeness in \mathcal{G} concerning discrepancy function $l_{f'}$ if for any $(f, g) \in \mathcal{H} \times (\mathcal{H} \cup \mathcal{G})$, it holds that

$$l_{f'}(f, g, \zeta) - l_{f'}(f, \mathcal{P}(f), \zeta) = \mathbb{E}_\zeta[l_{f'}(f, g, \zeta)], \quad (3.2)$$

where the expectation is taken over trajectory ζ . Besides, the operator satisfies that $\mathcal{P}(f^*) = f^*$.

The completeness assumption is an extension of Bellman completeness under the value-based hypothesis (Jin et al., 2021), incorporating the notion of the decomposition loss function (DLF) property proposed in Chen et al. (2022b). Our assumption diverges from the one posited in Zhong et al. (2022), where an auxiliary function class $\mathcal{P}(\mathcal{H}) \subseteq \mathcal{G}$ is introduced to enrich choices, accompanied with modifications tailored to accommodate the nuances of the average-reward setting.

Example 1 (Bellman completeness \subseteq Generalized completeness). Let the discrepancy function be the Bellman discrepancy in (3.1) with $\zeta_t = \{s_t, a_t, r_t, s_{t+1}\}$, and takes the (hypothesis-scheme) Bellman operator, defined as $\mathcal{T}(f) = \{\mathcal{T}_{J_f}(Q_f), J_f\}$ for all $f \in \mathcal{H}$, modified from (2.2). Then,

$$\begin{aligned} l_{f'}(f, g, \zeta_t) - l_{f'}(f, \mathcal{T}(f), \zeta_t) &= Q_g(s_t, a_t) - \mathcal{T}_{J_f}Q_f(s_t, a_t) \\ &= Q_g(s_t, a_t) - r(s_t, a_t) + \mathbb{E}_{\zeta_t}[V_f(s_{t+1})] + J_f = \mathbb{E}_{\zeta_t}[l_{f'}(f, g, \zeta_t)], \end{aligned}$$

where the expectation is taken over the transition state s_{t+1} from $\mathbb{P}(\cdot|s_t, a_t)$.

The preceding example illustrates that the Bellman discrepancy, a frequently employed discrepancy function across problems, satisfies both assumptions. More examples and choices of the discrepancy function for MLE-based algorithms are respectively provided in Appendix C and D.

3.2 Relation with Tractable Complexity Metric

To bridge the gap between concrete function approximation instances and the relatively abstract measure AGECE, this section introduces two intermediate metrics: the Eluder dimension (Russo and Van Roy, 2013) and the Average-reward Bellman Eluder (ABE) dimension. In particular, Chen et al. (2022a) employs the Eluder dimension to gauge the complexity of model-based hypothesis classes for infinite-horizon learning. To provide an intuitive complexity of value-based hypothesis classes, we additionally propose the ABE dimension, which is a generalization of the standard BE dimension (Jin et al., 2021). These two metrics provide valuable insights into the nature of AGECE.

Eluder Dimension We start with ϵ -independence notation (Russo and Van Roy, 2013).

Definition 4 (Point-wise ϵ -independence). Let \mathcal{H} be a function class defined on \mathcal{X} and consider sequence $\{z, x_1, \dots, x_n\} \in \mathcal{X}$. We say z is ϵ -independent of $\{x_1, \dots, x_n\}$ with respect to \mathcal{H} if there exists $f, f' \in \mathcal{H}$ such that $\sqrt{\sum_{i=1}^n (f(x_i) - f'(x_i))^2} \leq \epsilon$, but $|f(z) - f'(z)| \geq \epsilon$.

Based on ϵ -independence, the Eluder dimension can be efficiently defined as below.

Definition 5 (Eluder dimension). Let \mathcal{H} be a function class defined on \mathcal{X} . The Eluder dimension $\text{dim}_E(\mathcal{H}, \epsilon)$ is the length of the longest sequence $\{x_1, \dots, x_n\} \subset \mathcal{X}$ such that there exists $\epsilon' \geq \epsilon$ where x_i is ϵ' -independent of $\{x_1, \dots, x_{i-1}\}$ for all $i \in [n]$.

The following lemma shows that a model-based hypothesis class with a low Eluder dimension has low AGECE. Motivated by Ayoub et al. (2020); Chen et al. (2022a), we consider the Eluder dimension over function class derived from the model-based hypotheses class \mathcal{H} , defined as

$$\mathcal{X}_{\mathcal{H}} := \{X_{f,f'}(s, a) = (r_f + \mathbb{P}_{f'}V_f)(s, a) : f, f' \in \mathcal{H}\}.$$

Note that [Chen et al. \(2022a\)](#) considered function class $\mathcal{X}_{\mathcal{H},\mathcal{V}} := \{\mathbb{P}_f V(s, a) : f \in \mathbb{P}(\mathcal{H}), V \in \mathcal{V}\}$, where $\mathbb{P}(\mathcal{H})$ denotes the hypotheses class over the transition kernel and \mathcal{V} denotes the hypotheses class over the optimal bias function. We remark that definitions over function class based on (\mathbb{P}_f, V) with $(f, V) \in \mathcal{H} \times \mathcal{V}$ and (\mathbb{P}_f, r_f) with $f \in \mathcal{H}$ is almost equivalent and in this paper we focus on $\mathcal{X}_{\mathcal{H}}$ under the latter framework, aligning with the model-based hypothesis (see [Definition 2](#)).

Lemma 3.1 (Low Eluder dim \subseteq Low AGECE). Consider the discrepancy function, defined as

$$l_{f'}(f, g, \zeta_t) = (r_g + \mathbb{P}_g V_{f'})(s_t, a_t) - r(s_t, a_t) + V_{f'}(s_{t+1}), \quad (3.3)$$

and the expectation is taken over the transition state s_{t+1} from $\mathbb{P}(\cdot|s_t, a_t)$. Let $d_E = \dim_E(\mathcal{X}_{\mathcal{H}}, \epsilon)$ be the ϵ -Eluder dimension defined over $\mathcal{X}_{\mathcal{H}}$, then we have $d_G \leq 2d_E \cdot \log T$ and $\kappa_G \leq d_E$.

Average-Reward Bellman Eluder (ABE) Dimension Before delving into details of the average-reward BE (ABE) dimension, we start with two useful notations, distributional ϵ -independence and distributional Eluder (DE) dimension proposed by [Jin et al. \(2021\)](#), which is a generalization of point-wise ϵ -independence and Eluder dimension defined above (see [Definitions 4](#) and [5](#)).

Definition 6 (Distributional ϵ -independence). Let \mathcal{H} be a function class defined on \mathcal{X} and sequence $\{v, \mu_1, \dots, \mu_n\}$ be the probability measures over \mathcal{X} . We say v is ϵ -independent of $\{\mu_1, \dots, \mu_n\}$ with respect to \mathcal{H} if there exists $f \in \mathcal{H}$ such that $\sqrt{\sum_{i=1}^n (\mathbb{E}_{\mu_i}[f])^2} \leq \epsilon$, but $|\mathbb{E}_v[f]| \geq \epsilon$.

Definition 7 (Distributional Eluder dimension). Let \mathcal{H} be a function class defined on \mathcal{X} and Γ be a family of probability measures over \mathcal{X} . The distributional Eluder dimension $\dim_{DE}(\mathcal{H}, \Gamma, \epsilon)$ is the length of the longest sequence $\{\rho_1, \dots, \rho_n\} \subset \Gamma$ such that there exists $\epsilon' \geq \epsilon$ where ρ_i is ϵ' -independent of the remaining distribution sequence $\{\rho_1, \dots, \rho_{i-1}\}$ for all $i \in [n]$.

Now we are ready to introduce the average-reward Bellman Eluder (ABE) dimension. It is defined as the distributional Eluder (DE) dimension of average-reward Bellman error in [\(2.3\)](#).

Definition 8 (ABE dimension). Denote $\mathcal{E}_{\mathcal{H}} = \{\mathcal{E}(f)(s, a) : f \in \mathcal{H}\}$ be the collection of average-reward Bellman errors defined over $\mathcal{S} \times \mathcal{A}$. For any constant $\epsilon > 0$, the ϵ -ABE dimension of given hypotheses class \mathcal{H} is defined as $\dim_{ABE}(\mathcal{H}, \epsilon) := \dim_{DE}(\mathcal{E}_{\mathcal{H}}, \mathcal{D}_{\Delta}, \epsilon)$.

The lemma below posits that the value-based hypothesis problem with a low ABE dimension shall have a low AGECE in terms of the Bellman discrepancy.

Lemma 3.2 (Low ABE dim \subseteq Low AGECE). Consider the Bellman discrepancy function as defined in [\(3.1\)](#), and the expectation is taken over the transition state s_{t+1} from $\mathbb{P}(\cdot|s_t, a_t)$. Let $d_{ABE} = \dim_{ABE}(\mathcal{H}, \epsilon)$, then we have $d_G \leq 2d_{ABE} \cdot \log T$ and $\kappa_{ABE} \leq d_{ABE}$.

The Eluder dimension and ABE dimension can capture numerous concrete problems, respectively under model-based and value-based scenarios. Specifically, the Eluder dimension incorporates rich model-based problems like linear mixture AMDPs ([Wu et al., 2022](#)), and the ABE dimension can characterize tabular AMDPs ([Jaksch et al., 2010](#)), linear AMDPs ([Wei et al., 2021](#)), AMDPs with Bellman Completeness, generalized linear AMDPs, and kernel AMDPs, where the latter three problems are newly proposed for AMDPs. Details about the concrete examples are deferred to [Appendix D](#). Combining these facts and [Lemmas 3.1](#) and [3.2](#), we can conclude that AGECE serves as a unified complexity measure, as it encompasses all of these tractable AMDPs illustrated above.

4 Local-fitted Optimization with Optimism

To solve AMDPs with low AGECE value (see Definition 3), we propose the algorithm Local-fitted Optimization with Optimism (LOOP), whose pseudocode is given in Algorithm 1. At a high level, LOOP is a modified version of the classical fitted Q-iteration (Szepesvári, 2010) with optimistic planning and lazy policy updates. That is, the policies are only updated when a certain criterion is met (Line 2). When this is the case, LOOP performs three main steps:

- **Optimistic planning** (Line 4.1): Compute the most optimistic $f_t \in \mathcal{H}$ within \mathcal{B}_t that maximizes the corresponding average-reward J_t by solving a constrained optimization problem.
- **Construct confidence set** (Line 4.2): Construct the confidence set \mathcal{B}_t for optimization using \mathcal{D}_{t-1} , where all $f_t \in \mathcal{H}$ satisfying $\mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g) \leq \beta$ is included. Here, $\beta > 0$ defines the radius, corresponding to the log covering number of the hypothesis class \mathcal{H} .
- **Execute Policy and Update** Υ_t (Line 8-10): Choose the greedy policy $\pi_t = \pi_{f_t}$ as the exploration policy. Execute policy, collect data, and update trigger $\Upsilon_t = \mathcal{L}_{\mathcal{D}_t}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_t}(f_t, g)$.

Note that both the confidence set \mathcal{B}_t and the update condition Υ_t are constructed upon the (cumulative) squared discrepancy, which is crucial to the algorithmic design. It takes the form

$$\mathcal{L}_{\mathcal{D}_t}(f, f) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_t}(f, g), \quad \text{where } \mathcal{L}_{\mathcal{D}_t}(f, g) = \sum_{(f', \zeta) \in \mathcal{D}_t} \|l_{f'}(f, g, \zeta)\|_2^2, \quad (4.1)$$

where f' and ζ are drawn from \mathcal{D}_t , and $l_{f'}(f, g, \zeta)$ represents the discrepancy function, which varies across different RL problems. The cumulative squared discrepancy serves as an empirical estimation of the in-sample training error (see Definition 3). Besides, we highlight two key designs of LOOP:

- **Consistent control over discrepancy**: The construction of confidence set \mathcal{B}_t ensures that the cumulative squared discrepancy is controlled at level β in each step. To see this, suppose $\tau_t = t$, i.e., policy switches at the t -th step, then (4.2) ensures that $\mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g) \leq \beta$. Otherwise, if we do not switch policy at step t , then we must have $\Upsilon_{t-1} = \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g) \leq 4\beta$, as hypothesis $f_t = f_{t-1}$ remains unchanged.
- **Lazy policy update**: The regret decomposition in (5.1) elucidates that each policy switch incurs an additional cost of $|V_{t+1}(s_{t+1}) - V_t(s_{t+1})|$ in regret at each step (see (5.3)). This underscores the necessity of implementing lazy updates to achieve sublinear regret. Within the LOOP framework, policy updates occur adaptively, triggered only when a substantial increase in cumulative discrepancy surpassing 3β has occurred since the last update. Intuitively, a policy switch occurs when there is a notable infusion of new information from the recently collected data. Importantly, such gap is pivotal as it provides the theoretical foundation for the implementation of lazy updates, leveraging the problem's transferability structure (see (ii), Definition 3). LOOP employs a threshold of 4β , considering inherent uncertainty and estimation errors between the minimizer g and $\mathcal{P}(f_t)$, ensuring that the out-of-sample error will exceed β under the updating rule.

Similar to previous works in general function approximation (Jin et al., 2021; Du et al., 2021), our algorithm lacks a computationally efficient solution for constrained optimization problems. Instead, our focus is on the sample efficiency, as guaranteed by the theorem below.

Algorithm 1 Local-fitted Optimization with Optimism - LOOP($\mathcal{H}, \mathcal{G}, T, \delta$)

Parameter: Initial s_1 , span $\text{sp}(V^*)$, optimistic parameter $\beta = c \log(T\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(1/T)/\delta) \cdot \text{sp}(V^*)$

Initialize: Draw $a_1 \sim \text{Unif}(\mathcal{A})$ and set $\tau_0 \leftarrow 0$, $\Upsilon_0 \leftarrow 0$, $\mathcal{B}_0 \leftarrow \emptyset$, $\mathcal{D}_0 \leftarrow \emptyset$.

- 1: **for** $t = 1, \dots, T$ **do**
- 2: **if** $t = 1$ or $\Upsilon_{t-1} \geq 4\beta$ **then**
- 3: Set $\tau_t = t$.
- 4: Solve optimization problem $f_t = \arg\max_{f \in \mathcal{B}_t} J_{f_t}$, where

$$\mathcal{B}_t = \left\{ f \in \mathcal{H} : \mathcal{L}_{\mathcal{D}_{t-1}}(f, f) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f, g) \leq \beta \right\}, \quad (4.2)$$

- 5: Compute $Q_t = Q_{f_t}$, $V_t = V_{f_t}$ and $J_t = J_{f_t}$.
 - 6: **else**
 - 7: Retain $(f_t, J_t, V_t, Q_t, \tau_t) = (f_{t-1}, J_{t-1}, V_{t-1}, Q_{t-1}, \tau_{t-1})$.
 - 8: Take $a_t = \arg\max_{a \in \mathcal{A}} Q_t(s_t, a)$.
 - 9: Observe $r_t = r(s_t, a_t)$ and transition state s_{t+1} .
 - 10: Update $\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{(s_t, a_t, r_t, s_{t+1}; f_t)\}$ and $\Upsilon_t = \mathcal{L}_{\mathcal{D}_t}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_t}(f_t, g)$.
-

Theorem 4.1 (Regret). Under Assumptions 1-4, there exists constant c such that for any $\delta \in (0, 1)$ and horizon T , with probability at least $1 - 5\delta$, the regret of LOOP satisfies that

$$\text{Reg}(T) \leq \mathcal{O}\left(\text{sp}(V^*) \cdot d\sqrt{T\beta}\right),$$

where $\beta = c \log(T\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(1/T)/\delta) \cdot \text{sp}(V^*)$ and $d = \max\{\sqrt{d_G}, \kappa_G\}$. Here, $(d_G, \kappa_G) = \text{AGEC}(\mathcal{H}, \{l_f\}, 1/\sqrt{T})$ are AGECE defined in Definition 3, \mathcal{G} is the auxiliary function class defined in Definition 4, and $\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}(\cdot)$ denotes the covering number as defined in Definition 16.

Theorem 4.1 asserts that both value-based and model-based problems with low AGECE are tractable. Our algorithm LOOP achieves a $\tilde{\mathcal{O}}(\sqrt{T})$ regret and the multiplicative factor depends on span $\text{sp}(V^*)$, problem complexity $\max\{\sqrt{d_G}, \kappa_G\}$ and the log covering number. The proof sketch is provided in Section 5 and the detailed proof is deferred to Appendix E.

5 Proof Overview of Regret Analysis

In this section, we present the proof sketch of Theorem 4.1. In Section 4, we elucidated the construction of the confidence set and the circumstances in which updates are performed. Here, we delve into the theoretical analysis to substantiate the necessity of such designs.

Optimism and Regret Decomposition In LOOP, we apply an optimization based method to ensure the optimism $J_t \geq J^*$ at each step $t \in [T]$. Based on the optimistic algorithm, we propose a new regret decomposition method motivated by the standard performance difference lemma (PDL)

in episodic setting (Jiang et al., 2017), following the form as below:

$$\text{Reg}(T) \leq \underbrace{\sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t)}_{\text{Bellman error}} + \underbrace{\sum_{t=1}^T \mathbb{E}_{s_{t+1} \sim \mathbb{P}(\cdot|s_t, a_t)}[V_t(s_{t+1})] - V_t(s_t)}_{\text{Realization error}}. \quad (5.1)$$

Bound over Bellman Error The control over the Bellman error is achieved through the design of a confidence set and update condition that combinely controls the empirical squared discrepancy. Note that the construction of the confidence set filters f_t with a limited sum of empirical squared discrepancy. Note that $\mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g)$ can be regarded as an empirical overestimation of the squared discrepancy, controlled at $\mathcal{O}(\beta)$ regardless of updating. Then,

$$\text{In-sample training error} = \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \lesssim \beta \quad \forall t \in [T], \quad (5.2)$$

with high probability, and β is pre-determined optimistic parameter depends on horizon T and the log ρ -covering number. Recall that the dominance coefficient d_G regulates that Bellman error $\lesssim \left[d_G \sum_{t=1}^T \left(\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \right) \right]^{1/2}$, thus we have Bellman error $\leq \mathcal{O}(\sqrt{d_G \beta T})$.

Bound over Realization error The realization error is small if the switching cost is low as the concentration arguments indicated that with high probability it holds

$$\text{Realization error} \leq \underbrace{\text{sp}(V^*) \cdot \mathcal{N}(T)}_{\text{Switching cost}} + \underbrace{\mathcal{O}(\text{sp}(V^*) \cdot \sqrt{T \log(1/\delta)})}_{\text{Azuma-Hoeffding term}}, \quad (5.3)$$

where $\mathcal{N}(t)$ denote Switching cost defined as $\mathcal{N}(T) = \#\{t \in [T] : \tau_t \neq \tau_{t-1}\}$. Motivated by the recent work of Xiong et al. (2023), the main idea of low-switching control is summarized below. The key step is that the minimizer g is a good approximator of $\mathcal{P}(f_t)$ such that

$$0 \leq \mathcal{L}_{\mathcal{D}_t}(f_t, \mathcal{P}(f_t)) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_t}(f_t, g) \leq \beta, \quad \forall t \in [T] \quad (5.4)$$

with high probability based on the minimization and definition of optimistic parameter β . In the following analysis, we assume that (5.4) holds true. Consider an update occurs at step $t+1$, it implies that $\mathcal{L}_{\mathcal{D}_t}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_t}(f_t, g) > 4\beta$ and the latest update at step τ_t ensures that $\mathcal{L}_{\mathcal{D}_{\tau_t-1}}(f_{\tau_t}, f_{\tau_t}) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{\tau_t-1}}(f_{\tau_t}, g) \leq \beta$. Combine (5.4) with arguments above, we have

$$(i). \mathcal{L}_{\mathcal{D}_t}(f_t, f_t) - \mathcal{L}_{\mathcal{D}_t}(f_t, \mathcal{P}(f_t)) > 3\beta, \quad (ii). \mathcal{L}_{\mathcal{D}_{\tau_t-1}}(f_{\tau_t}, f_{\tau_t}) - \mathcal{L}_{\mathcal{D}_{\tau_t-1}}(f_{\tau_t}, \mathcal{P}(f_{\tau_t})) \leq \beta. \quad (5.5)$$

Using the concentration argument and (i), (ii) in (5.5), the out-sample training error between updates is lower bounded by $\sum_{i=\tau_t}^t \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_i, f_i, \zeta_i)]\|_2^2 > \beta$. Let $b_1, \dots, b_{\mathcal{N}(T)+1}$ be the sequence of updated steps, take summation over the whole process and we have

$$\mathcal{N}(T)\beta \leq \sum_{u=1}^{\mathcal{N}(T)} \sum_{t=b_u}^{b_{u+1}-1} \|\mathbb{E}_{\zeta_t}[l_{f_t}(f_t, f_t, \zeta_t)]\|_2^2 = \sum_{t=1}^T \|\mathbb{E}_{\zeta_t}[l_{f_t}(f_t, f_t, \zeta_t)]\|_2^2 \leq \mathcal{O}(\kappa_G \cdot \beta \log T), \quad (5.6)$$

where the first inequality follows arguments above and the second is based on the definition of transferability (see Definition 3) given $\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \leq \mathcal{O}(\beta)$ for all $t \in [T]$. Thus, we have $\mathcal{N}(T) \leq \mathcal{O}(\kappa_G \log T)$ and the Realization error is bounded by $\mathcal{O}(\kappa_G \cdot \text{sp}(V^*) \log T)$. Please refer to Lemma E.3 in Appendix E.4 for a formal statement and detailed techniques.

6 Conclusion

This work studies the infinite-horizon average-reward MDPs under general function approximation. To address the unique challenges of AMDPs, we introduce a new complexity metric — average-reward generalized eluder coefficient (AGEC) and a unified algorithm named Local-fitted Optimization with OPTimism (LOOP). We posit that our work paves the way for future work, including developing more general frameworks for AMDPs and new algorithms with sharper regret bounds.

References

- Abbasi-Yadkori, Y., Bartlett, P., Bhatia, K., Lazic, N., Szepesvari, C., and Weisz, G. (2019). PoliteX: Regret bounds for policy iteration using expert prediction. In *International Conference on Machine Learning*, pages 3692–3702. PMLR.
- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. (2011). Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24.
- Agarwal, A., Hsu, D., Kale, S., Langford, J., Li, L., and Schapire, R. (2014). Taming the monster: A fast and simple algorithm for contextual bandits. In *International Conference on Machine Learning*, pages 1638–1646. PMLR.
- Agarwal, A., Jin, Y., and Zhang, T. (2022). Vo q l: Towards optimal regret in model-free rl with nonlinear function approximation. *arXiv preprint arXiv:2212.06069*.
- Auer, P., Jaksch, T., and Ortner, R. (2008). Near-optimal regret bounds for reinforcement learning. *Advances in neural information processing systems*, 21.
- Ayoub, A., Jia, Z., Szepesvari, C., Wang, M., and Yang, L. (2020). Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pages 463–474. PMLR.
- Bai, Y., Xie, T., Jiang, N., and Wang, Y.-X. (2019). Provably efficient q-learning with low switching cost. *Advances in Neural Information Processing Systems*, 32.
- Bartlett, P. L. and Tewari, A. (2012). Regal: A regularization based algorithm for reinforcement learning in weakly communicating mdps. *arXiv preprint arXiv:1205.2661*.
- Cai, Q., Yang, Z., Jin, C., and Wang, Z. (2020). Provably efficient exploration in policy optimization. In *International Conference on Machine Learning*, pages 1283–1294. PMLR.
- Chen, X., Hu, J., Jin, C., Li, L., and Wang, L. (2022a). Understanding domain randomization for sim-to-real transfer. *International Conference on Learning Representations*.

- Chen, Z., Li, C. J., Yuan, A., Gu, Q., and Jordan, M. I. (2022b). A general framework for sample-efficient function approximation in reinforcement learning. *arXiv preprint arXiv:2209.15634*.
- Dani, V., Hayes, T. P., and Kakade, S. M. (2008). Stochastic linear optimization under bandit feedback.
- Dann, C., Mohri, M., Zhang, T., and Zimmert, J. (2021). A provably efficient model-free posterior sampling method for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 34:12040–12051.
- Domingues, O. D., Ménard, P., Pirotta, M., Kaufmann, E., and Valko, M. (2021). A kernel-based approach to non-stationary reinforcement learning in metric spaces. In *International Conference on Artificial Intelligence and Statistics*, pages 3538–3546. PMLR.
- Du, S., Kakade, S., Lee, J., Lovett, S., Mahajan, G., Sun, W., and Wang, R. (2021). Bilinear classes: A structural framework for provable generalization in rl. In *International Conference on Machine Learning*, pages 2826–2836. PMLR.
- Foster, D. J., Golowich, N., and Han, Y. (2023). Tight guarantees for interactive decision making with the decision-estimation coefficient. *arXiv preprint arXiv:2301.08215*.
- Foster, D. J., Kakade, S. M., Qian, J., and Rakhlin, A. (2021). The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*.
- Fruit, R., Pirotta, M., Lazaric, A., and Ortner, R. (2018). Efficient bias-span-constrained exploration-exploitation in reinforcement learning. In *International Conference on Machine Learning*, pages 1578–1586. PMLR.
- Hao, B., Lazic, N., Abbasi-Yadkori, Y., Joulani, P., and Szepesvári, C. (2021). Adaptive approximate policy iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 523–531. PMLR.
- He, J., Zhao, H., Zhou, D., and Gu, Q. (2022). Nearly minimax optimal reinforcement learning for linear markov decision processes. *arXiv preprint arXiv:2212.06132*.
- Hernández-Lerma, O. (2012). *Adaptive Markov control processes*, volume 79. Springer Science & Business Media.
- Hu, J., Zhong, H., Jin, C., and Wang, L. (2022). Provable sim-to-real transfer in continuous domain with partial observations. *arXiv preprint arXiv:2210.15598*.
- Huang, J., Zhong, H., Wang, L., and Yang, L. F. (2023). Tackling heavy-tailed rewards in reinforcement learning with function approximation: Minimax optimal and instance-dependent regret bounds. *arXiv preprint arXiv:2306.06836*.
- Jaksch, T., Ortner, R., and Auer, P. (2010). Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(51):1563–1600.

- Jiang, N., Krishnamurthy, A., Agarwal, A., Langford, J., and Schapire, R. E. (2017). Contextual decision processes with low bellman rank are pac-learnable. In *International Conference on Machine Learning*, pages 1704–1713. PMLR.
- Jin, C., Liu, Q., and Miryoosefi, S. (2021). Bellman eluder dimension: New rich classes of rl problems, and sample-efficient algorithms. *Advances in neural information processing systems*, 34:13406–13418.
- Jin, C., Liu, Q., and Yu, T. (2022). The power of exploiter: Provable multi-agent rl in large state spaces. In *International Conference on Machine Learning*, pages 10251–10279. PMLR.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR.
- Kong, D., Salakhutdinov, R., Wang, R., and Yang, L. F. (2021). Online sub-sampling for reinforcement learning with general function approximation. *arXiv preprint arXiv:2106.07203*.
- Krishnamurthy, A., Agarwal, A., and Langford, J. (2016). Pac reinforcement learning with rich observations. *Advances in Neural Information Processing Systems*, 29.
- Lazic, N., Yin, D., Abbasi-Yadkori, Y., and Szepesvari, C. (2021). Improved regret bound and experience replay in regularized policy iteration. In *International Conference on Machine Learning*, pages 6032–6042. PMLR.
- Li, X. and Sun, Q. (2023). Variance-aware robust reinforcement learning with linear function approximation with heavy-tailed rewards. *arXiv preprint arXiv:2303.05606*.
- Liu, Q., Chung, A., Szepesvári, C., and Jin, C. (2022). When is partially observable reinforcement learning not scary? In *Conference on Learning Theory*, pages 5175–5220. PMLR.
- Liu, Q., Netrapalli, P., Szepesvari, C., and Jin, C. (2023a). Optimistic mle: A generic model-based algorithm for partially observable sequential decision making. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, pages 363–376.
- Liu, Z., Lu, M., Xiong, W., Zhong, H., Hu, H., Zhang, S., Zheng, S., Yang, Z., and Wang, Z. (2023b). One objective to rule them all: A maximization objective fusing estimation and planning for exploration. *arXiv preprint arXiv:2305.18258*.
- Ortner, R. (2020). Regret bounds for reinforcement learning via markov chain concentration. *Journal of Artificial Intelligence Research*, 67:115–128.
- Ouyang, Y., Gagrani, M., Nayyar, A., and Jain, R. (2017). Learning unknown markov decision processes: A thompson sampling approach. *Advances in neural information processing systems*, 30.
- Proper, S. and Tadepalli, P. (2006). Scaling model-based average-reward reinforcement learning for product delivery. In *European Conference on Machine Learning*, pages 735–742. Springer.
- Puterman, M. L. (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.

- Russo, D. and Van Roy, B. (2013). Eluder dimension and the sample complexity of optimistic exploration. *Advances in Neural Information Processing Systems*, 26.
- Srinivas, N., Krause, A., Kakade, S. M., and Seeger, M. (2009). Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*.
- Sun, W., Jiang, N., Krishnamurthy, A., Agarwal, A., and Langford, J. (2019). Model-based rl in contextual decision processes: Pac bounds and exponential improvements over model-free approaches. In *Conference on learning theory*, pages 2898–2933. PMLR.
- Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Szepesvári, C. (2010). Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103.
- Wang, J., Wang, M., and Yang, L. F. (2022). Near sample-optimal reduction-based policy learning for average reward mdp. *arXiv preprint arXiv:2212.00603*.
- Wang, R., Salakhutdinov, R. R., and Yang, L. (2020). Reinforcement learning with general value function approximation: Provably efficient approach via bounded eluder dimension. *Advances in Neural Information Processing Systems*, 33:6123–6135.
- Wang, Y., Wang, R., Du, S. S., and Krishnamurthy, A. (2019). Optimism in reinforcement learning with generalized linear function approximation. *arXiv preprint arXiv:1912.04136*.
- Wei, C.-Y., Jahromi, M. J., Luo, H., and Jain, R. (2021). Learning infinite-horizon average-reward mdps with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 3007–3015. PMLR.
- Wei, C.-Y., Jahromi, M. J., Luo, H., Sharma, H., and Jain, R. (2020). Model-free reinforcement learning in infinite-horizon average-reward markov decision processes. In *International conference on machine learning*, pages 10170–10180. PMLR.
- Wu, Y., Zhou, D., and Gu, Q. (2022). Nearly minimax optimal regret for learning infinite-horizon average-reward mdps with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 3883–3913. PMLR.
- Xiong, N., Yang, Z., and Wang, Z. (2023). A general framework for sequential decision-making under adaptivity constraints. *arXiv preprint arXiv:2306.14468*.
- Yang, L. and Wang, M. (2019). Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004. PMLR.
- Zanette, A., Lazaric, A., Kochenderfer, M., and Brunskill, E. (2020). Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pages 10978–10989. PMLR.
- Zhang, Z. and Ji, X. (2019). Regret minimization for reinforcement learning by evaluating the optimal bias function. *Advances in Neural Information Processing Systems*, 32.

- Zhang, Z. and Xie, Q. (2023). Sharper model-free reinforcement learning for average-reward markov decision processes. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 5476–5477. PMLR.
- Zhang, Z., Zhou, Y., and Ji, X. (2020). Almost optimal model-free reinforcement learning via reference-advantage decomposition. *Advances in Neural Information Processing Systems*, 33:15198–15207.
- Zhao, H., He, J., Zhou, D., Zhang, T., and Gu, Q. (2023). Variance-dependent regret bounds for linear bandits and reinforcement learning: Adaptivity and computational efficiency. *arXiv preprint arXiv:2302.10371*.
- Zhong, H., Xiong, W., Zheng, S., Wang, L., Wang, Z., Yang, Z., and Zhang, T. (2022). Gec: A unified framework for interactive decision making in mdp, pomdp, and beyond. *arXiv preprint arXiv:2211.01962*.
- Zhong, H. and Zhang, T. (2023). A theoretical analysis of optimistic proximal policy optimization in linear markov decision processes. *arXiv preprint arXiv:2305.08841*.
- Zhou, D. and Gu, Q. (2022). Computationally efficient horizon-free reinforcement learning for linear mixture mdps. *arXiv preprint arXiv:2205.11507*.
- Zhou, D., Gu, Q., and Szepesvari, C. (2021a). Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*, pages 4532–4576. PMLR.
- Zhou, D., He, J., and Gu, Q. (2021b). Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*, pages 12793–12802. PMLR.

A Backgrounds and Technical Novelties

Infinite-horizon Average-reward MDPs. Pioneering works by [Auer et al. \(2008\)](#) and [Bartlett and Tewari \(2012\)](#) laid the foundation for model-based algorithms operating within the online framework with sub-linear regret. In recent years, the pursuit of improved regret guarantees has led to the emergence of a multitude of new algorithms. In tabular case, these advancements include numerous model-based approaches ([Ouyang et al., 2017](#); [Fruit et al., 2018](#); [Zhang and Ji, 2019](#); [Ortner, 2020](#)) and model-free algorithms ([Abbasi-Yadkori et al., 2019](#); [Wei et al., 2020](#); [Hao et al., 2021](#); [Lazic et al., 2021](#); [Zhang and Xie, 2023](#)). In the context of function approximation, POLITEX ([Abbasi-Yadkori et al., 2019](#)), a variant of the regularized policy iteration, is the first model-free algorithm with linear value-function approximation, and achieves $\tilde{O}(T^{\frac{3}{4}})$ regret for the ergodic MDP. The work by [Hao et al. \(2021\)](#) followed the same setting and improved the results to $\tilde{O}(T^{\frac{2}{3}})$ with an adaptive approximate policy iteration (AAPI) algorithm. [Wei et al. \(2021\)](#) proposed an optimistic Q-learning algorithm FOPO for the linear function approximation, and achieve a near-optimal $\tilde{O}(\sqrt{T})$ regret. On another line of research, [Wu et al. \(2022\)](#) delved into the linear function approximation under the framework of linear mixture model, which is mutually uncoverable concerning linear MDPs ([Wei et al., 2021](#)), and proposed UCRL2-VTR based on the value-targeted regression ([Ayoub et al., 2020](#)). Recent work of [Chen et al. \(2022a\)](#) expanded the scope of research by addressing the general function approximation problem in average-reward RL and proposed the SIM-TO-REAL algorithm, which can be regarded as an extension to UCRL2-VTR. In comparison to the works mentioned, our algorithm, LOOP, not only addresses all the problems examined in those studies but also extends its applicability to newly identified models. See Table 1 for a summary.

Function Approximation in Finite-horizon MDPs. In the pursuit of developing sample-efficient algorithms capable of handling large state spaces, extensive research efforts have converged on the linear function approximation problems within the finite-horizon setting. See [Yang and Wang \(2019\)](#); [Wang et al. \(2019\)](#); [Jin et al. \(2020\)](#); [Ayoub et al. \(2020\)](#); [Cai et al. \(2020\)](#); [Zhou et al. \(2021a,b\)](#); [Zhou and Gu \(2022\)](#); [Agarwal et al. \(2022\)](#); [He et al. \(2022\)](#); [Zhong and Zhang \(2023\)](#); [Zhao et al. \(2023\)](#); [Huang et al. \(2023\)](#); [Li and Sun \(2023\)](#) and references therein. Furthermore, [Wang et al. \(2020\)](#) studied RL with general function approximation and adopted the eluder dimension ([Russo and Van Roy, 2013](#)) as a complexity measure. Before this, [Jiang et al. \(2017\)](#) considered a substantial subset of problems with low Bellman ranks. Building upon these foundations, [Jin et al. \(2021\)](#) combined both the Eluder dimension and Bellman error, thereby broadening the scope of solvable problems under the concept of the Bellman Eluder (BE) dimension. In a parallel line of research, [Sun et al. \(2019\)](#) proposed the witness ranking focusing on the low-rank structures, and [Du et al. \(2021\)](#) extended it to encompass more scenarios with the bilinear class. Besides, [Foster et al. \(2021, 2023\)](#) provided a unified framework, decision estimation coefficient, for interactive decision making. The work of [Chen et al. \(2022b\)](#) extended the value-based GOLF ([Jin et al., 2021](#)) with the introduction of the discrepancy loss function to handle the broader admissible Bellman characterization (ABC) class. More recently, [Zhong et al. \(2022\)](#); [Liu et al. \(2023b\)](#) proposed a unified framework measured by generalized eluder coefficient (GEC), an extension to [Dann et al. \(2021\)](#) that captures almost all known tractable problems. All these works are restricted to the finite-horizon regime, and their complexity measure and algorithms are not applicable in the infinite-

horizon average-reward setting.

Low-Switching Cost Algorithms. Addressing low-switching cost problems in bandit and reinforcement learning has seen notable progress. Abbasi-Yadkori et al. (2011) first proposed an algorithm for linear bandits with $\mathcal{O}(\log T)$ switching cost. Subsequent research extended this to tabular MDPs, including works of Bai et al. (2019); Zhang et al. (2020). A significant stride was made by Kong et al. (2021), who introduced importance scores to handle low-switching cost scenarios in general function approximation with complexity measured by eluder dimension (Russo and Van Roy, 2013). Recently, Xiong et al. (2023) introduced the eluder condition (EC) class, offering a comprehensive framework to address all tractable low-switching cost problems above. In the context of average-reward RL, Wei et al. (2021); Wu et al. (2022); Chen et al. (2022a); Hu et al. (2022) developed low-switching algorithms to control the regret under linear structure or model-based class, leaving a unifying framework for both value-based and model-based problems an open problem.

Further Elaboration on Our Contributions and Technical Novelties. Compared to episodic MDPs or discounted MDPs, AMDPs present unique challenges that prevent a straightforward extension of existing algorithms and analyses from these well-studied domains. One notable distinction is a different regret notion in average-reward RL due to a different form of the Bellman optimality equation. Furthermore, such a difference is coupled with the challenge of exploration in the context of general function approximation. To effectively bound this regret, we introduce a new regret decomposition approach within the context of general function approximation (refer to (5.1) and (5.3)). This regret decomposition suggests that the total regret can be controlled by the cumulative Bellman error and the switching cost. Inspired by this, we propose an optimistic algorithm with lazy updates in the general function approximation setting, which uses the residue of the loss function as the indicator for deciding when to conduct policy updates. Such a lazy policy update scheme adaptively divides the total of T steps into $\mathcal{O}(\log T)$ epochs, which is significantly different from (OLSVI.FH; Wei et al., 2021) that reduces the infinite-horizon setting to the finite-horizon setting by splitting the whole learning procedure into several H -length epoch, where H typically chosen as $\Theta(\sqrt{T})$ (Wei et al., 2021). We remark that such an adaptive lazy updating design and corresponding analysis are pivotal in achieving the optimal $\tilde{\mathcal{O}}(\sqrt{T})$ rate, as opposed to the $\tilde{\mathcal{O}}(T^{3/4})$ regret in (OLSVI.FH; Wei et al., 2021). Moreover, our approach is an extension to the existing lazy update approaches for average-reward setting (Wei et al., 2021; Wu et al., 2022) that leverages the postulated linear structure and is not applicable to problems with general function approximation. Furthermore, to accommodate the average-reward term, we introduce a new complexity measure AGECE, which characterizes the exploration challenge in general function approximation. Compared with Zhong et al. (2022), our additional transferability restriction is tailored for the infinite-horizon setting and plays a crucial role in analyzing the low-switching error. Despite this additional transferability restriction, AGECE can still serve as a unifying complexity measure in the infinite-horizon average-reward setting, like the role of GEC in the finite-horizon setting. Specifically, AGECE captures a rich class of tractable AMDP models, including all previously recognized AMDPs, including all known tractable AMDPs, and some newly identified AMDPs. See Table 1 for a summary.

Algorithm 2 MLE-based Local-fitted Optimization with Optimism - MLE-LOOP(\mathcal{H}, T, δ)

Parameter: Initial s_1 , span $\text{sp}(V^*)$, optimistic parameter $\beta = c \log(T\mathcal{N}_{\mathcal{H}}(1/T)/\delta)$

Initialize: Draw $a_1 \sim \text{Unif}(\mathcal{A})$ and set $\tau_0 \leftarrow 0$, $\Upsilon_0 \leftarrow 0$, $\mathcal{B}_0 \leftarrow \emptyset$, $\mathcal{D}_0 \leftarrow \emptyset$.

- 1: **for** $t = 1, \dots, T$ **do**
- 2: **if** $t = 1$ or $\Upsilon_{t-1} \geq 3\sqrt{\beta t}$ **then**
- 3: Update $\tau_t = t$.
- 4: Solve optimization problem $f_t = \arg\max_{f \in \mathcal{B}_t} J_{f_t}$, where

$$\mathcal{B}_t = \left\{ f \in \mathcal{H} : \mathcal{L}_{\mathcal{D}_{t-1}}(f, f) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f, g) \leq \beta \right\}, \quad (\text{C.1})$$

- 5: Update $Q_t = Q_{f_t}$, $V_t = V_{f_t}$, $J_t = J_{f_t}$ and $g_t = \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g)$.
 - 6: **else**
 - 7: Retain $(f_t, g_t, J_t, V_t, Q_t, \tau_t) = (f_{t-1}, g_{t-1}, J_{t-1}, V_{t-1}, Q_{t-1}, \tau_{t-1})$.
 - 8: Take $a_t = \arg\max_{a \in \mathcal{A}} Q_t(s_t, a)$.
 - 9: Collect reward $r_t = r(s_t, a_t)$ and transition state s_{t+1} .
 - 10: Update $\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{(s_t, a_t, r_t, s_{t+1})\}$, $\Upsilon_t = \sum_{(s,a) \in \mathcal{D}_t} \text{TV}(\mathbb{P}_{f_t}(\cdot|s, a), \mathbb{P}_{g_t}(\cdot|s, a))$.
-

B Discussions

Discussion about distribution families Beyond the singleton distribution family \mathcal{D}_{Δ} taken in this paper, there exists a notable distribution family $\mathcal{D}_{\mathcal{H}} = \{\mathcal{D}_{\mathcal{H},t}\}_{t \in [T]}$, proposed in Jin et al. (2021), where $\mathcal{D}_{\mathcal{H},t}$ characterizes probability measure over $\mathcal{S} \times \mathcal{A}$ obtained by executing different policies induced by $f_1, \dots, f_{t-1} \in \mathcal{H}$, measures the detailed distribution under sequential policies. However, in this paper, we exclude the consideration of $\mathcal{D}_{\mathcal{H}}$ for two principal reasons. First, evaluations of average-reward RL focus on the difference between *observed* rewards $r(s_t, a_t)$ and optimal average reward J^* — as opposed to the expected value V_h^π (i.e. *expected* sum of reward) under specific policy and optimal value at step $h \in [H]$ in episodic setting — rendering the introduction of $\mathcal{D}_{\mathcal{H}}$ unnecessary. Second, in infinite settings, the measure of such distribution becomes highly intricate and impractical given *different* policy induced by f_1, \dots, f_T over a potentially infinite T -steps. As a comparison, in the episode setting a *fixed* policy induced by f_t is executed over a finite H -step.

Discussion about V-type variant Previous research has demonstrated the existence of extensive classes of MDPs characterized by a low V-type BE dimension but cannot be captured by a infinite Q-type BE dimension in episodic setting (Jiang et al., 2017; Jin et al., 2021), and similar scenarios can be straightforwardly extended to AMDPs. However, in this paper, we abstain from introducing the V-type variant, because to the best of our knowledge, solving a V-type problem relies on the *off-policy strategy* or auxiliary simulators, which is unfortunately infeasible under the infinite-horizon online learning procedure. We maintain this for potential solutions in future research endeavors.

C Alternative Choices of Discrepancy Function

Note that there is another line of research that addresses model-based problems using Maximum Likelihood Estimator (MLE)-based approaches (Liu et al., 2023a; Zhong et al., 2022), as opposed to the value-targeted regression. We remark that these MLE-based approaches can be also incorporated within our framework through the use of the discrepancy function:

$$l_{f'}(f, g, \zeta_t) = \frac{1}{2} |\mathbb{P}_f(s_{t+1}|s_t, a_t) / \mathbb{P}_{f^*}(s_{t+1}|s_t, a_t) - 1|, \quad (\text{C.2})$$

where the trajectory is $\zeta_t = (s_t, a_t, s_{t+1})$ with expectation taken over the next transition state s_{t+1} from $\mathbb{P}(\cdot|s_t, a_t)$ such that $\mathbb{E}_{\zeta_t}[l_{f'}(f, g, \zeta_t)] = \text{TV}(\mathbb{P}_f(\cdot|s_t, a_t), \mathbb{P}_{f^*}(\cdot|s_t, a_t))$. To accommodate the discrepancy function in (C.2), we introduce a natural variant of AGECE defined below.

Definition 9 (MLE-AGECE). Given hypothesis class \mathcal{H} , discrepancy function $\{l_f\}_{f \in \mathcal{H}}$ in (C.2) and constant $\epsilon > 0$, the average-reward generalized eluder coefficients MLE-AGECE($\mathcal{H}, \{l_f\}, \epsilon$) is defined as the smallest coefficients κ_G and d_G , such that following two conditions hold with absolute constants $C_1, C_2 > 0$:

- (i) (Bellman dominance) There exists constant $d_G > 0$ such that

$$\sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t) \leq \left[d_G \cdot \sum_{t=1}^T \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \right]^{1/2} + C_1 \cdot \text{sp}(V^*) \min\{d_G, T\} + T\epsilon.$$

- (ii) (MLE-Transferability) There exists constant $\kappa_G > 0$ such that for hypotheses $f_1, \dots, f_T \in \mathcal{H}$, if it holds that $\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \leq \beta$ for all $t \in [T]$, then we have

$$\sum_{t=1}^T \|\mathbb{E}_{\zeta_t}[l_{f_t}(f_t, f_t, \zeta_t)]\|_1 \leq \text{poly}(\log T) \sqrt{\kappa_G \cdot \beta T} + C_2 \cdot \text{sp}(V^*)^2 \min\{\kappa_G, T\} + 2T\epsilon^2.$$

The main difference between the MLE-based variant (see Definition 9) and the original AGECE (see Definition 3) is that the transferability coefficient is defined over the l_1 -norm of out-sample training error rather than the l_2 -norm, and similar condition is considered in Liu et al. (2023a); Xiong et al. (2023). Now we are ready to introduce the algorithm for the alternative discrepancy function in (C.2) and please see Algorithm 2 for complete pseudocode. The main modification lies in the construction of confidence set and update condition Υ_t . Here, the confidence set follows

$$\mathcal{B}_t = \{f_t \in \mathcal{H} : \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g) \leq \beta\}, \quad \mathcal{L}_{\mathcal{D}}(f, g) = - \sum_{(s, a, s') \in \mathcal{D}} \log \mathbb{P}_g(s'|s, a).$$

In comparison, the update condition follows that $\Upsilon_t = \sum_{(s, a) \in \mathcal{D}_t} \text{TV}(\mathbb{P}_{f_t}(\cdot|s, a), \mathbb{P}_{g_t}(\cdot|s, a))$. Unlike the standard LOOP algorithm, the the confidence set and update condition in the MLE-based variant no longer shares the same construction. Here, we explicitly redesign the algorithm for MLE-based approaches, and the theoretical guarantee is provided below.

Theorem C.1 (Cumulative regret). Under Assumptions 1-2 and the discrepancy function in (C.2) with self-completeness such that $\mathcal{G} = \mathcal{H}$, there exists constant c such that for any $\delta \in (0, 1)$ and time horizon T , with probability at least $1 - 4\delta$, the regret of MLE-LOOP satisfies that

$$\text{Reg}(T) \leq \mathcal{O}\left(\text{sp}(V^*) \cdot d\sqrt{T\beta}\right),$$

where $\beta = c \log(T\mathcal{B}_{\mathcal{H}}(1/T)/\delta)$, $d = \max\{\sqrt{d_G}, \kappa_G\}$. $(d_G, \kappa_G) = \text{MLE-AGEC}(\mathcal{H}, \{l_f\}, 1/\sqrt{T})$ denotes MLE-AGEC defined in Definition 9 and $\mathcal{B}_{\mathcal{H}}(\cdot)$ denotes the bracketing number.

The proof of Theorem C.1 is similar to that of Theorem 4.1, and can be found in Appendix I.1.

D Concrete Examples

In this section, we present concrete examples of problems for AMDP. We remark that the understanding of function approximation problems under the average-reward setting is quite limited, and to our best knowledge, existing works have primarily focused on linear approximation (Wei et al., 2021; Wu et al., 2022) and model-based general function approximation (Chen et al., 2022a). Here, we introduce a variety of function classes with low AGEK. Beyond the examples considered in existing work, these newly proposed function classes are mostly natural extensions from their counterpart the finite-horizon episode setting (Jin et al., 2020; Zanette et al., 2020; Du et al., 2021; Domingues et al., 2021), which can be extended to the average-reward problems with moderate justifications.

D.1 Linear Function Approximation and Variants

Linear function approximation Consider the linear FA, which encompasses a wide range of concrete problems with state-action bias function linear in a d -dimensional feature mapping. Specifically, the linear function class \mathcal{H}_{Lin} is formally defined as $\mathcal{H}_{\text{Lin}} = \{Q(\cdot, \cdot) = \langle \omega, \phi(\cdot, \cdot) \rangle, J \in \mathbb{J}(\mathcal{H}) \mid \|\omega\|_2 \leq \frac{1}{2}\text{sp}(V^*)\sqrt{d}\}$, where $\|\phi(s, a)\|_2 \leq \sqrt{2}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and first coordinate fixed to 1. We emphasize that such scaling is without loss of generality justified in Lemma H.8. To begin with, we introduce two problems: linear AMDP and AMDP with linear Bellman completion.

Definition 10 (Linear AMDP, Wei et al. (2021)). There exists a known feature mapping $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$, an unknown d -dimensional signed measures $\mu = (\mu_1, \dots, \mu_d)$ over \mathcal{S} , and an unknown reward parameter $\theta \in \mathbb{R}^d$, such that the transition kernel the reward function can be written as

$$\mathbb{P}(\cdot | s, a) = \langle \phi(s, a), \mu(\cdot) \rangle, \quad r(s, a) = \langle \phi(s, a), \theta \rangle. \quad (\text{D.1})$$

for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Without loss of generality, we assume that the feature mapping satisfies that $\|\phi(s, a)\|_2 \leq \sqrt{2}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ with $\phi_{(1)} = 1$ fixed, $\|\theta\|_2 \leq \sqrt{d}$ and $\|\mu(\mathcal{S})\|_2 \leq \sqrt{d}$, where denote $\mu(\mathcal{S}) = (\mu_1(\mathcal{S}), \dots, \mu_d(\mathcal{S}))$ and $\mu_i(\mathcal{S}) = \int_{\mathcal{S}} d\mu_i(s)$ be the total measure of \mathcal{S} .

We remark that the scaling on the feature mapping can help in overcoming the gap between the episodic setting and the average-reward one without incurring any additional computational or

statistical costs. To illustrate the necessity, note that for linear AMDP we have

$$\begin{aligned} Q^*(s, a) &= r(s, a) + \mathbb{E}_{s' \sim \mathbb{P}(s, a)}[V^*(s')] - J^* \quad (\text{see (2.1)}) \\ &= \phi(s, a)^\top \left(\theta - J^* \mathbf{e}_{(1)} + \int_{\mathcal{S}} V^*(s') d\mu(s') \right) := \langle \phi(s, a), \omega \rangle, \end{aligned} \quad (\text{D.2})$$

where denote $\mathbf{e}_{(1)} = (1, 0, \dots, 0) \in \mathbb{R}^d$. Next, we provide the AMDPs with linear Bellman completion, modified from [Zanette et al. \(2020\)](#), which is a more general setting than linear AMDPs.

Definition 11 (Linear Bellman completion). There exists a known feature mapping $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ such that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $J \in \mathbb{J}(\mathcal{H})$, we have

$$\langle \mathcal{T}_J(\omega), \phi(s, a) \rangle := r(s, a) + \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} \left[\max_{a' \in \mathcal{A}} \left\{ \omega^\top \phi(s', a') \right\} \right] - J, \quad (\text{D.3})$$

which indicates that the function remains linear in feature mapping under the Bellman operator.

Generalized linear function approximation To introduce the nonlinearity beyond linear FA, we expand the linear function class \mathcal{H}_{Lin} by incorporating a link function. In the context of generalized linear FA, the function class is defined as $\mathcal{H}_{\text{Glin}} = \{Q(\cdot, \cdot) = \sigma(\omega^\top \phi(\cdot, \cdot)), J \in \mathbb{J}(\mathcal{H}) \mid \|\omega\|_2 \leq \sqrt{d}\}$, where $\|\phi(s, a)\|_2 \leq 1$ holds for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, and $\sigma : \mathbb{R} \mapsto \mathbb{R}$ is an α -bi-Lipschitz link function such that $|\sigma(x)| \leq \frac{1}{2}\text{sp}(V^*)$ for all $x \in \mathbb{R}$. Formally, σ is α -bi-Lipschitz continuous if

$$\frac{1}{\alpha}|x - x'| \leq |\sigma(x) - \sigma(x')| \leq \alpha|x - x'|, \quad \forall x, x' \in \mathbb{R}. \quad (\text{D.4})$$

We remark that the generalized linear function class $\mathcal{H}_{\text{Glin}}$ degenerates to the standard linear function class \mathcal{H}_{Lin} if we choose $\sigma(x) = x$. Modified from [Wang et al. \(2019\)](#) for the episodic setting, we define AMDPs with generalized linear Bellman completion as follows.

Definition 12 (Generalized linear Bellman completion). There exists a known feature mapping $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ such that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $J \in \mathbb{J}(\mathcal{H})$, we have

$$\sigma \left(\mathcal{T}_J(\omega)^\top \phi(\cdot, \cdot) \right) := r(s, a) + \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} \left[\max_{a' \in \mathcal{A}} \left\{ \sigma \left(\omega^\top \phi(s', a') \right) \right\} \right] - J, \quad (\text{D.5})$$

where σ is an α -bi-Lipschitz function with $|\sigma(x)| \leq \frac{1}{2}\text{sp}(V^*)$ for all $x \in \mathbb{R}$, and $\alpha \geq 1$.

The proposition below states that (generalized) linear function classes have low AGECE.

Proposition D.1 (Linear FA \subset Low AGECE). Consider linear function class \mathcal{H}_{Lin} and generalized linear function class $\mathcal{H}_{\text{Glin}}$ with a d -dimensional feature mapping $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$, if the problem follows one of Definitions 10-12, then it have low AGECE under Bellman discrepancy in (3.1):

$$d_G \leq \mathcal{O}(d \log(\text{sp}(V^*)\sqrt{d}\epsilon^{-1}) \log T), \quad \kappa_G \leq \mathcal{O}(d \log(\text{sp}(V^*)\sqrt{d}\epsilon^{-1})).$$

Linear Q^*/V^* AMDP Moreover, we consider the linear Q^*/V^* AMDPs, which is modified from the one in [Du et al. \(2021\)](#) under the episodic setting. Note that the linear Q^*/V^* AMDPs are

not strictly a value-based problem (see Definition 1), as there exists an additional assumption made beyond (Q^*, J^*) to impose the linear structure of state bias function V^* .

Definition 13 (Linear Q^*/V^* AMDP). There exists known feature mappings $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^{d_1}$, $\psi : \mathcal{S} \mapsto \mathbb{R}^{d_2}$, and unknown vectors $\omega^* \in \mathbb{R}^{d_1}$, $\theta^* \in \mathbb{R}^{d_2}$ such that optimal value functions follow

$$Q^*(s, a) = \langle \phi(s, a), \omega^* \rangle, \quad V^*(s') = \langle \psi(s'), \theta^* \rangle,$$

for all $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$. Without loss of generality, we assume the feature mappings hold $\|\phi(s, a)\|_2 \leq \sqrt{2}$, $\|\psi(s, a)\|_2 \leq \sqrt{2}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ with first coordinate $\phi_{(1)} = \psi_{(1)} = 1$ fixed, and the parameters follow $\|\omega^*\|_2 \leq \frac{1}{2}\text{sp}(V^*)\sqrt{d_1}$ and $\|\theta^*\|_2 \leq \frac{1}{2}\text{sp}(V^*)\sqrt{d_2}$.

Note that for linear Q^*/V^* AMDPs, if we disregard the information related to the feature mapping $\psi : \mathcal{S} \mapsto \mathbb{R}^{d_2}$, the problem is a standard linear FA problem (with further assumption required). To fully leverage the structural information provided in the linear Q^*/V^* structure, AGECC can also offer an appropriate complexity measure for this variant of the value-based problem, as stated below.

Proposition D.2 (Linear $Q^*/V^* \subset \text{Low AGECC}$). Linear Q^*/V^* AMDPs with coupled (d_1, d_2) -dimensional feature mappings $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^{d_1}$ and $\psi : \mathcal{S} \mapsto \mathbb{R}^{d_2}$ have low AGECC such that

$$d_G \leq \mathcal{O}(\underline{d} \log(\text{sp}(V^*)\sqrt{\underline{d}}\epsilon^{-1}) \log T), \quad \kappa_G \leq \mathcal{O}(\underline{d} \log(\text{sp}(V^*)\sqrt{\underline{d}}\epsilon^{-1})),$$

where denote $\underline{d} = d_1 + d_2$ be the sum of dimensions of both feature mappings.

The proposition above asserts that in linear Q^*/V^* AMDPs, acquiring additional structural information through the state bias function comes at an added computational cost of $\tilde{\mathcal{O}}(d_2)$ in comparison to standard linear FA. Notably, linear FA, generalized linear FA, and linear Q^*/V^* AMDPs are representative value-based problems. The proof of this proposition relies on the ABE dimension as an intermediate complexity measure, and further result directly follows Lemma 3.2.

D.2 Kernel Function Approximation

In this subsection, we first introduce the notion of effective dimension. With this notion, we prove a useful proposition that any linear kernel function class with a low effective dimension also has low AGECC. Consider the kernel FA, a natural extension to linear FA (see Appendix D.1 for detailed), from the d -dimensional Euclidean space \mathbb{R}^d to a separable kernel Hilbert space \mathcal{K} . Note that the linear kernel function class is defined as $\mathcal{H}_{\text{Ker}} = \{Q(\cdot, \cdot) = \langle \phi(\cdot, \cdot), \omega \rangle_{\mathcal{K}}, J \in \mathbb{J}(\mathcal{H}) \mid \|\omega\|_{\mathcal{K}} \leq \text{sp}(V^*)R\}$ given a separable Hilbert space \mathcal{K} , where the feature mapping $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{K}$ satisfies that $\|\phi(s, a)\|_{\mathcal{K}} \leq 1$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. To better measure the complexity of problems in Hilbert space \mathcal{K} with potentially infinite dimension, we introduce the ϵ -effective dimension.

Definition 14 (ϵ -effective dimension). Consider a set \mathcal{Z} with the possibly infinite elements in a given separable Hilbert space \mathcal{K} , the ϵ -effective dimension, denoted by $\text{dim}_{\text{eff}}(\mathcal{Z}, \epsilon)$, is defined as the length n of the longest sequence satisfying the condition below:

$$\sup_{z_1, \dots, z_n \in \mathcal{Z}} \left\{ \frac{1}{n} \log \det \left(\mathbf{I} + \frac{1}{\epsilon^2} \sum_{t=1}^n z_t z_t^\top \right) \leq e^{-1} \right\}.$$

Here, the concept of ϵ -effective dimension is inspired by the measurement of maximum information gain (Srinivas et al., 2009) and is later introduced as a complexity measure of Hilbert space in Du et al. (2021); Zhong et al. (2022). Similar to Jin et al. (2021), we strengthen the assumption over completeness and require the \mathcal{H} to be self-complete concerning the operator such that $\mathcal{G} = \mathcal{H}$ to simplify the analysis of the complexity of \mathcal{H}_{Ker} . We're now ready to show kernel FA has low AGEK.

Proposition D.3 (Kernel FA \subset Low AGEK). Under the self-completeness, kernel FA with function class \mathcal{H}_{Ker} concerning a known feature mapping $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{K}$ have low AGEK such that

$$d_G \leq \dim_{\text{eff}}(\mathcal{X}, \epsilon/2\text{sp}(V^*)R) \log T, \quad \kappa_G \leq \dim_{\text{eff}}(\mathcal{X}, \epsilon/2\text{sp}(V^*)R),$$

where denote $\mathcal{X} = \{\phi(s, a) : (s, a) \in \mathcal{S} \times \mathcal{A}\}$ as the collection of feature mappings.

The proposition above demonstrates that the kernel FA with a low ϵ -effective dimension over the Hilbert space has low AGEK. As a special case of kernel FA, if we choose $\mathcal{K} = \mathbb{R}^d$, then we can prove that the RHS in the proposition above is upper bounded by $\tilde{\mathcal{O}}(d)$ (Jin et al., 2022).

D.3 Linear Mixture AMDP

In this subsection, we focus on the average-reward linear mixture problem considered in Wu et al. (2022). In this context, the hypotheses function class is defined as $\mathcal{H}_{\text{LM}} = \{\mathbb{P}(s'|s, a) = \theta^\top \phi(s, a, s'), r(s, a) = \theta^\top \psi(s, a) \mid \|\theta\|_2 \leq 1\}$ with known feature mappings $\phi : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}^d$, $\psi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$, and an unknown parameter $\theta \in \mathbb{R}^d$. Detailedly, the problem is defined as below.

Definition 15 (Linear mixture AMDPs, Wu et al. (2022)). There exists a known feature mapping $\phi : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}^d$, $\psi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$, and unknown vectors $\theta \in \mathbb{R}^d$, it holds that

$$\mathbb{P}(s'|s, a) = \langle \theta, \phi(s, a, s') \rangle, \quad r(s, a) = \langle \theta, \psi(s, a) \rangle,$$

for all $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$. Without loss of generality, we assume that the feature mappings satisfy that $\|\phi(s, a)\|_2 \leq \sqrt{d}$ and $\|\psi(s, a)\|_2 \leq \sqrt{d}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Now we show that the linear mixture problem is tractable under the framework of AGEK.

Proposition D.4 (Linear mixture \subset Low AGEK). Consider linear mixture problem with hypotheses class \mathcal{H}_{Lin} and d -dimensional feature (ϕ, ψ) . If the discrepancy function is chosen as

$$l_{f'}(f, g, \zeta_t) = \theta_g^\top \left\{ \psi(s_t, a_t) + \int_{\mathcal{S}} \phi(s_t, a_t, s') V_{f'}(s') ds' \right\} - r(s_t, a_t) - V_{f'}(s_{t+1}), \quad (\text{D.6})$$

and takes $\mathcal{H} = \mathcal{G}$ with operator following $\mathcal{P}(f) = f^*$ for all $f \in \mathcal{H}$, it has low AGEK such that

$$d_G \leq \mathcal{O} \left(d \log \left(\text{sp}(V^*)T/\sqrt{d\epsilon} \right) \right), \quad \kappa_G \leq \mathcal{O} \left(d \log \left(\text{sp}(V^*)T/\sqrt{d\epsilon} \right) \right).$$

The proposition posits that AGEK can capture the linear mixture AMDP, based on a modified version of the Bellman discrepancy function in (2.2). In contrast to the linear FA discussed in

Appendix D.1 above, the presence of the average-reward term in this model-based problem does not impose any additional computational or statistical burden, and there is no need for structural assumptions on feature mappings, such as a fixed first coordinate, considering discrepancy in (D.6).

E Proof of Main Results for LOOP

E.1 Proof of Theorem 4.1

Proof of Theorem 4.1. Note that the regret can be decomposed as

$$\begin{aligned}
\text{Reg}(T) &= \sum_{t=1}^T \left(J^* - r(s_t, a_t) \right) \leq \sum_{t=1}^T \left(J_t - r(s_t, a_t) \right) && \text{(optimism)} \\
&\stackrel{(a)}{=} \sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t) - \sum_{t=1}^T \mathbb{E}_{s_{t+1} \sim \mathbb{P}(\cdot | s_t, a_t)} \left[Q_t(s_t, a_t) - \max_{a \in \mathcal{A}} Q_t(s_{t+1}, a) \right] \\
&\stackrel{(b)}{=} \underbrace{\sum_{i=1}^T \mathcal{E}(f_t)(s_t, a_t)}_{\text{Bellman error}} + \underbrace{\sum_{t=1}^T \left[\mathbb{E}_{s_{t+1} \sim \mathbb{P}(\cdot | s_t, a_t)} [V_t(s_{t+1})] - V_t(s_t) \right]}_{\text{Realization error}}, && \text{(E.1)}
\end{aligned}$$

where step (a) and step (b) follow the definition of the Bellman optimality operator and the greedy policy. Below, we will present the bound of Bellman error and Realization error respectively.

Step 1: Bound over Bellman error Recall that the of confidence set ensures that $\mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g) \leq \mathcal{O}(\beta)$ across all steps. Using the concentration arguments, we can infer

$$\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i} [l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \leq \mathcal{O}(\beta), \quad \text{(E.2)}$$

with high probability and the formal statements are deferred to Lemma E.2 in Appendix E.3. In the following arguments, we assume the above event holds. Take $\epsilon = 1/\sqrt{T}$, recall the definition of dominance coefficient d_G in AGE $C(\mathcal{H}, \mathcal{J}, l, \epsilon)$ and it directly indicates that

$$\text{Bellman error} = \sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t) \leq \left[d_G \sum_{t=1}^T \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i} [l_{f_i}(f_t, f_t, \zeta)]\|_2^2 \right]^{1/2} + \mathcal{O} \left(\text{sp}(V^*) \sqrt{d_G T} \right),$$

and thus the Bellman error can be upper bounded by $\mathcal{O} \left(\text{sp}(V^*) \sqrt{d_G \beta T} \right)$.

Step 2: Bound over Realization error To bound the Realization error, we shall use the concentration argument and the upper bound of the switching cost. Note that

$$\begin{aligned}
\text{Realization error} &\stackrel{(c)}{=} \sum_{t=1}^T [V_t(s_{t+1}) - V_t(s_t)] + \mathcal{O}(\text{sp}(V^*)\sqrt{T \log(1/\delta)}), \\
&= \sum_{t=1}^T [V_{\tau_t}(s_{t+1}) - V_{\tau_{t+1}}(s_{t+1})] + \mathcal{O}(\text{sp}(V^*)\sqrt{T \log(1/\delta)}), \quad (\text{Shift}) \\
&= \sum_{t=1}^T [V_{\tau_t}(s_{t+1}) - V_{\tau_{t+1}}(s_{t+1})] \mathbb{1}(\tau_t \neq \tau_{t+1}) + \mathcal{O}(\text{sp}(V^*)\sqrt{T \log(1/\delta)}) \\
&\leq \text{sp}(V^*) \cdot \mathcal{N}(T) + \mathcal{O}(\text{sp}(V^*)\sqrt{T \log(1/\delta)}) \stackrel{(d)}{\leq} \mathcal{O}(\text{sp}(V^*) \cdot \kappa_G \sqrt{T \log(1/\delta)}), \quad (\text{E.3})
\end{aligned}$$

where step (c) directly follows the Azuma-Hoeffding inequality and step (d) is based the fact that $\|V_{\tau_t} - V_{\tau_{t+1}}\|_\infty \leq \text{sp}(V^*)$ and the bounded switching cost such that $\mathcal{N}(T) \leq \mathcal{O}(\kappa_G \log T)$, where κ_G is the transferability coefficient in AGEK with $\epsilon = 1/\sqrt{T}$. Please refer to Lemma E.3 in Appendix E.4 for the detailed statement and proof of the bounded switching cost.

Step 3: Combine the bounded erroes Plugging (E.2) and (E.3) back into n (E.1), we have

$$\begin{aligned}
\text{Reg}(T) &\leq \text{Bellman error} + \text{Realization error} \\
&\leq \mathcal{O}(\text{sp}(V^*)\sqrt{d_G \beta T}) + \mathcal{O}(\text{sp}(V^*)\kappa_G \sqrt{T \log(1/\delta)}) = \mathcal{O}(\text{sp}(V^*) \cdot d \sqrt{T \beta}),
\end{aligned}$$

where denote $d = \max\{\sqrt{d_G}, \kappa_G\}$ is a function of $(d_G, \kappa_G) = \text{AGEK}(\mathcal{H}, \{l_f\}_{f \in \mathcal{H}}, 1/\sqrt{T})$. In the arguments above, the optimistic parameter is chosen as $\beta = c \log(T \mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(1/T)/\delta) \cdot \text{sp}(V^*)$, which takes the upper bound of the optimistic parameters, aligning with the choice in Lemma E.1, Lemma E.2, and Lemma E.3. Then finish the proof of cumulative regret for LOOP in Algorithm 1. \square

E.2 Proof of Lemma E.1

Lemma E.1 (Optimism). Under Assumptions 1-4, LOOP is an optimistic algorithm such that it ensures $J_t \geq J^*$ for all $t \in [T]$ with probability greater than $1 - \delta$.

Proof of Lemma E.1. Denote $\mathcal{V}_\rho(\mathcal{G})$ be the ρ -cover of \mathcal{G} and $\mathcal{N}_\rho(\mathcal{G})$ be the size of ρ -cover $\mathcal{V}_\rho(\mathcal{G})$. Consider fixed $(i, g) \in [T] \times \mathcal{G}$ and define the auxiliary function

$$X_{i, f_i}(g) := \|l_{f_i}(f^*, g, \zeta_i)\|_2^2 - \|l_{f_i}(f^*, f^*, \zeta_i)\|_2^2, \quad (\text{E.4})$$

where f^* is the optimal hypothesis in value-based problems and the true hypothesis in model-based ones. Let \mathcal{F}_t be the filtration induced by $\{s_1, a_1, \dots, s_t, a_t\}$ and note that f_1, \dots, f_t is fixed under

the filtration, then we have

$$\begin{aligned}
\mathbb{E}[X_{i,f_i}(g)|\mathcal{F}_i] &= \mathbb{E}_{\zeta_i}[\|l_{f_i}(f^*, g, \zeta_i)\|_2^2 - \|l_{f_i}(f^*, \mathcal{P}(f^*), \zeta_i)\|_2^2 | \mathcal{F}_i] \\
&= \mathbb{E}_{\zeta_i} \left[[l_{f_i}(f^*, g, \zeta_i) - l_{f_i}(f^*, \mathcal{P}(f^*), \zeta_i)] \cdot [l_{f_i}(f^*, g, \zeta_i) + l_{f_i}(f^*, \mathcal{P}(f^*), \zeta_i)] \middle| \mathcal{F}_i \right] \\
&= \mathbb{E}_{\zeta_i} \left[\mathbb{E}_{\zeta_i} [l_{f_i}(f^*, g, \zeta_i)] \cdot [l_{f_i}(f^*, g, \zeta_i) + l_{f_i}(f^*, \mathcal{P}(f^*), \zeta_i)] \middle| \mathcal{F}_i \right] \\
&= \|\mathbb{E}_{\zeta_i} [l_{f_i}(f^*, g, \zeta_i)]\|_2^2,
\end{aligned}$$

where the equation follows the definition of generalized completeness (see Assumption 4):

$$\begin{cases} \mathbb{E}_{\zeta_i} [l_{f'}(f, g, \zeta)] = l_{f'}(f, g, \zeta) - l_{f'}(f, \mathcal{P}(f), \zeta), \\ \mathbb{E}_{\zeta_i} [l_{f'}(f, g, \zeta)] = \mathbb{E}_{\zeta_i} [l_{f'}(f, g, \zeta) + l_{f'}(f, \mathcal{P}(f), \zeta)]. \end{cases}$$

Similarly, we can obtain that the second moment of the auxiliary function is bounded by

$$\mathbb{E}[X_{i,f_i}(g)^2 | \mathcal{F}_i] \leq \mathcal{O}\left(sp(v^*)^2 \|\mathbb{E}_{\zeta_i} [l_{f_i}(f^*, g, \zeta_i)]\|_2^2\right),$$

By Freedman's inequality (see Lemma H.7), with probability greater than $1 - \delta$ it holds that

$$\begin{aligned}
&\left| \sum_{i=1}^t X_{i,f_i}(g) - \sum_{i=1}^t \|\mathbb{E}_{\zeta_i} [l_{f_i}(f^*, g, \zeta_i)]\|_2^2 \right| \\
&\leq \mathcal{O} \left(\sqrt{\log(1/\delta) \cdot sp(V^*)^2 \sum_{i=1}^t \|\mathbb{E}_{\zeta_i} [l_{f_i}(f^*, g, \zeta_i)]\|_2^2} + \log(1/\delta) \right).
\end{aligned}$$

By taking union bound over $[T] \times \mathcal{V}_\rho(\mathcal{G})$, for any $(t, \phi) \in [T] \times \mathcal{V}_\rho(\mathcal{G})$ we have $-\sum_{i=1}^t X_{i,f_i}(\phi) \leq \mathcal{O}(\zeta)$, where $\zeta = sp(V^*) \log(T\mathcal{N}_{\mathcal{G}}(\rho)/\delta)$ and we use the fact that $\|\mathbb{E}_{\zeta_i} [l_{f_i}(f^*, g, \zeta_i)]\|_2^2$ is non-negative. Recall the definition of ρ -cover, it ensures that for any $g \in \mathcal{G}$, there exists $\phi \in \mathcal{V}_\epsilon(\mathcal{G})$ such that $\|g(s, a) - \phi(s, a)\|_1 \leq \rho$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Therefore, for any $g \in \mathcal{G}$ we have

$$-\sum_{i=1}^t X_{i,f_i}(g) \leq \mathcal{O}(\zeta + t\rho). \quad (\text{E.5})$$

Combine the (E.5) above and the designed confidence set, then for all $t \in [T]$ it holds that

$$\mathcal{L}_{\mathcal{D}_{t-1}}(f^*, f^*) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f^*, g) = -\sum_{i=1}^{t-1} X_{i,f_i}(\tilde{g}) \leq \mathcal{O}(\zeta + t\rho) < \beta, \quad (\text{E.6})$$

where \tilde{g} is the local minimizer to $\mathcal{L}_{\mathcal{D}_{t-1}}(f^*, g)$, and we take the covering coefficient as $\rho = 1/T$ and optimistic parameter as $\beta = c \log(T\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(1/T)/\delta) sp(V^*)$. Based on (E.6), with probability greater than $1 - \delta$, f^* is a candidate of the confidence set such that $J_t \geq J^*$ for all $t \in [T]$. \square

E.3 Proof of Lemma E.2

Lemma E.2. For fixed $\rho > 0$ and the optimistic parameter $\beta = c(\text{sp}(V^*) \log(T\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(\rho)/\delta) + T\rho)$ where $c > 0$ is constant large enough, then it holds that

$$\sum_{i=1}^{t-1} \mathbb{E}_{\zeta_i} \|l_{f_i}(f_t, f_t, \zeta_i)\|^2 \leq \mathcal{O}(\beta), \quad (\text{E.7})$$

for all $t \in [T]$ with probability greater than $1 - \delta$.

Proof of Lemma E.2. Denote $\mathcal{V}_\rho(\mathcal{H})$ be the ρ -cover of \mathcal{H} and $\mathcal{N}_{\mathcal{H}}(\rho)$ be the size of ρ -cover $\mathcal{V}_\rho(\mathcal{H})$. Consider fixed $(i, f) \in [T] \times \mathcal{H}$ and define the auxiliary function

$$X_{i,f_i}(f) := \|l_{f_i}(f, f, \zeta_i)\|_2^2 - \|l_{f_i}(f, \mathcal{P}(f), \zeta_i)\|_2^2,$$

Let \mathcal{F}_t be the filtration induced by $\{s_1, a_1, \dots, s_t, a_t\}$ and note that f_1, \dots, f_t is fixed under the filtration, then we have

$$\begin{aligned} \mathbb{E}[X_{i,f_i}(f) | \mathcal{F}_i] &= \mathbb{E}_{\zeta_i} [\|l_{f_i}(f, f, \zeta_i)\|_2^2 - \|l_{f_i}(f, \mathcal{P}(f), \zeta_i)\|_2^2 | \mathcal{F}_i] \\ &= \mathbb{E}_{\zeta_i} \left[[l_{f_i}(f, f, \zeta_i) - l_{f_i}(f, \mathcal{P}(f), \zeta_i)] \cdot [l_{f_i}(f, f, \zeta_i) + l_{f_i}(f, \mathcal{P}(f), \zeta_i)] \middle| \mathcal{F}_i \right] \\ &= \mathbb{E}_{\zeta_i} [l_{f_i}(f, f, \zeta_i)] \cdot \mathbb{E}_{\zeta_i} [l_{f_i}(f, f, \zeta_i) + l_{f_i}(f, \mathcal{P}(f), \zeta_i) | \mathcal{F}_i] \\ &= \|\mathbb{E}_{\zeta_i} [l_{f_i}(f, f, \zeta_i)]\|_2^2, \end{aligned}$$

where the equation generalized completeness (see Lemma E.1). Similarly, we can obtain that the second moment of the auxiliary function is bounded by

$$\mathbb{E}[X_{i,f_i}(f)^2 | \mathcal{F}_i] \leq \mathcal{O}\left(\text{sp}(v^*)^2 \|\mathbb{E}_{\zeta_i} [l_{f_i}(f, f, \zeta_i)]\|_2^2\right),$$

By Freedman's inequality in Lemma H.7, with probability greater than $1 - \delta$ we have

$$\begin{aligned} &\left| \sum_{i=1}^t X_{i,f_i}(f) - \sum_{i=1}^t \|\mathbb{E}_{\zeta_i} [l_{f_i}(f, f, \zeta_i)]\|_2^2 \right| \\ &\leq \mathcal{O} \left(\sqrt{\log(1/\delta) \cdot \text{sp}(V^*)^2 \sum_{i=1}^t \|\mathbb{E}_{\zeta_i} [l_{f_i}(f, f, \zeta_i)]\|_2^2} + \log(1/\delta) \right). \end{aligned}$$

Define $\zeta = \text{sp}(V^*) \log(T\mathcal{N}_{\mathcal{H}}(\rho)/\delta)$, by taking a union bound over ρ -covering of hypothesis set \mathcal{H} , we can obtain that with probability greater than $1 - \delta$, for all $(t, \phi) \in [T] \times \mathcal{V}_\rho(\mathcal{H})$ we have

$$\begin{aligned} &\left| \sum_{i=1}^t X_{i,f_i}(\phi) - \sum_{i=1}^t \|\mathbb{E}_{\zeta_i} [l_{f_i}(\phi, \phi, \zeta_i)]\|_2^2 \right| \\ &\leq \mathcal{O} \left(\sqrt{\zeta \cdot \text{sp}(V^*)^2 \sum_{i=1}^t \|\mathbb{E}_{\zeta_i} [l_{f_i}(\phi, \phi, \zeta_i)]\|_2^2} + \zeta \right). \quad (\text{E.8}) \end{aligned}$$

The following analysis assumes that the event above is true. Recall that the LOOP ensures that

$$\begin{aligned}
\sum_{i=1}^{t-1} X_{i,f_i}(f_t) &= \sum_{i=1}^{t-1} \|l_{f_i}(f_t, f_t, \zeta_i)\|_2^2 - \sum_{i=1}^{t-1} \|l_{f_i}(f_t, \mathcal{P}(f_t), \zeta_i)\|_2^2, \\
&\leq \sum_{i=1}^{t-1} \|l_{f_i}(f_t, f_t, \zeta_i)\|_2^2 - \inf_{g \in \mathcal{G}} \sum_{i=1}^{t-1} \|l_{f_i}(f_t, g, \zeta_i)\|_2^2, \\
&= \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g) \leq \mathcal{O}(\beta),
\end{aligned} \tag{E.9}$$

where the last inequality is based on the confidence set and the update condition combined. Note that if the update is executed at time t , the confidence set ensures that

$$\mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g) \leq \beta,$$

within the update step t . Otherwise, if the update condition is not triggered, we have $f_{\tau_t} = f_t$ and

$$\Upsilon_{t-1} = \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{t-1}}(f_t, g) \leq 4\beta.$$

Recall that based on the definition of ρ -cover for any $f \in \mathcal{H}$, there exists $\phi \in \mathcal{V}_\rho(\mathcal{H})$ such that $\|g(s, a) - \phi(s, a)\|_1 \leq \rho$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have the in-sample training error is bounded by

$$\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \leq \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(\phi_t, \phi_t, \zeta_i)]\|_2^2 + \mathcal{O}(t\rho), \quad (\rho\text{-approximation})$$

$$= \sum_{i=1}^{t-1} X_{i,f_i}(\phi_t) + \mathcal{O}(t\rho + \zeta) \tag{E.8}$$

$$= \sum_{i=1}^{t-1} X_{i,f_i}(f_t) + \mathcal{O}(t\rho + \zeta) \leq \mathcal{O}(T\rho + \zeta + \beta) = \mathcal{O}(\beta), \tag{E.10}$$

where the last inequality follows (E.9), and takes $\beta = c((\text{sp}(V^*) \log(T\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(\rho)/\delta) + T\rho))$. \square

E.4 Proof of Lemma E.3

Lemma E.3. Let $\mathcal{N}(T)$ be the switching cost with time horizon T , defined as

$$\mathcal{N}(T) = \#\{t \in [T] : \tau_t \neq \tau_{t-1}\}.$$

Given fixed $\rho > 0$ and the optimistic parameter $\beta = c(\text{sp}(V^*) \log(T\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(\rho)/\delta) + T\rho)$, where $c > 0$ is large enough constant, then with probability greater than $1 - 2\delta$ we have

$$\mathcal{N}(T) \leq \mathcal{O}(\kappa_G \log T + \beta^{-1} T \log T \epsilon^2),$$

where κ_G is the transferability coefficient with respect to $\text{AGEC}(\mathcal{H}, \{l_{f'}\}, \epsilon)$.

Proof of Lemma E.3. Denote $\mathcal{V}_\rho(\mathcal{H})$ be the ρ -cover of \mathcal{H} and $\mathcal{N}_\mathcal{H}(\rho)$ be the size of ρ -cover $\mathcal{V}_\rho(\mathcal{H})$.

Step 1: Bound the difference of discrepancy between the minimizer and $\mathcal{P}(f)$.

Consider fixed tuple $(i, f, g) \in [T] \times \mathcal{H} \times \mathcal{G}$ and define auxiliary function as

$$X_{i,f_i}(f, g) := \|l_{f_i}(f, g, \zeta_i)\|_2^2 - \|l_{f_i}(f, \mathcal{P}(f), \zeta_i)\|_2^2$$

Let \mathcal{F}_t be the filtration induced by $\{s_1, a_1, \dots, s_t, a_t\}$ and note that f_1, \dots, f_t is fixed under the filtration, then we have

$$\begin{aligned} \mathbb{E}[X_{i,f_i}(f, g) | \mathcal{F}_i] &= \mathbb{E}_{\zeta_i}[\|l_{f_i}(f, g, \zeta_i)\|_2^2 - \|l_{f_i}(f, \mathcal{P}(f), \zeta_i)\|_2^2 | \mathcal{F}_i] \\ &= \mathbb{E}_{\zeta_i}\left[\left[l_{f_i}(f, g, \zeta_i) - l_{f_i}(f, \mathcal{P}(f), \zeta_i)\right] \cdot \left[l_{f_i}(f, g, \zeta_i) + l_{f_i}(f, \mathcal{P}(f), \zeta_i)\right] \middle| \mathcal{F}_i\right] \\ &= \mathbb{E}_{\zeta_i}[l_{f_i}(f, g, \zeta_i)] \cdot \mathbb{E}_{\zeta_i}[l_{f_i}(f, g, \zeta_i) + l_{f_i}(f, \mathcal{P}(f), \zeta_i) | \mathcal{F}_i] \\ &= \|\mathbb{E}_{\zeta_i}[l_{f_i}(f, g, \zeta_i)]\|_2^2, \end{aligned}$$

where the equation generalized completeness (see Lemma E.1). Similarly, we can obtain that the second moment of the auxiliary function is bounded by

$$\mathbb{E}[X_{i,f_i}(f, g)^2 | \mathcal{F}_i] \leq \mathcal{O}\left(\text{sp}(V^*)^2 \|\mathbb{E}_{\zeta_i}[l_{f_i}(f, g, \zeta_i)]\|_2^2\right),$$

By Freedman's inequality in Lemma H.7, with probability greater than $1 - \delta$

$$\begin{aligned} &\left| \sum_{i=1}^t X_{i,f_i}(f, g) - \sum_{i=1}^t \|\mathbb{E}_{\zeta_i}[l_{f_i}(f, g, \zeta_i)]\|_2^2 \right| \\ &\leq \mathcal{O}\left(\sqrt{\log(1/\delta) \cdot \text{sp}(V^*)^2 \sum_{i=1}^t \|\mathbb{E}_{\zeta_i}[l_{f_i}(f, g, \zeta_i)]\|_2^2} + \log(1/\delta)\right) \end{aligned}$$

Define $\zeta = \text{sp}(V^*) \log(T\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(\rho)/\delta)$, by taking a union bound over ρ -covering of hypothesis set $\mathcal{H} \times \mathcal{G}$, with probability greater than $1 - \delta$, for all $(t, \phi, \varphi) \in [T] \times \mathcal{V}_\rho(\mathcal{H}) \times \mathcal{V}_\rho(\mathcal{G})$ it holds

$$\begin{aligned} &\left| \sum_{i=1}^t X_{i,f_i}(\phi, \psi) - \sum_{i=1}^t \|\mathbb{E}_{\zeta_i}[l_{f_i}(\phi, \psi, \zeta_i)]\|_2^2 \right| \\ &\leq \mathcal{O}\left(\sqrt{\zeta \cdot \text{sp}(V^*)^2 \sum_{i=1}^t \|\mathbb{E}_{\zeta_i}[l_{f_i}(\phi, \psi, \zeta_i)]\|_2^2} + \zeta\right), \end{aligned} \tag{E.11}$$

where $\zeta = \text{sp}(V^*) \log(T\mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(\rho)/\delta)$. Note that $\|\mathbb{E}_{\zeta_i}[l_{f_i}(\phi, \psi, \zeta_i)]\|_2^2$ is non-negative, then it holds that $-\sum_{i=1}^t X_{i,f_i}(\phi, \varphi) \leq \mathcal{O}(\zeta)$ for all $t \in [T]$. Based on (E.11) and the ρ -approximation, we have

$$-\sum_{i=1}^t X_{i,f_i}(f, g) \leq \mathcal{O}(\zeta + t\rho), \quad \forall t \in [T],$$

for any $(f, g) \in \mathcal{H} \times \mathcal{G}$. Recall that $\beta = c \log(T \mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(\rho)/\delta) \text{sp}(V^*)$, for all $t \in [T]$ we have

$$\begin{aligned} \mathcal{L}_{\mathcal{D}_t}(f_t, \mathcal{P}(f_t)) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_t}(f_t, g) &= \sum_{i=1}^t \|l_{f_i}(f_t, \mathcal{P}(f_t), \zeta_i)\|_2^2 - \inf_{g \in \mathcal{G}} \sum_{i=1}^t \|l_{f_i}(f_t, g, \zeta_i)\|_2^2 \\ &= - \sum_{i=1}^t X_{i, f_i}(f_t, \tilde{g}) \leq \mathcal{O}(\zeta + t\rho) \leq \beta. \end{aligned} \quad (\text{E.12})$$

Combine (E.12), and the fact that g is defined as the local minimizer among auxiliary class \mathcal{G} and $\mathcal{P}(f_t) \in \mathcal{G}$, then for all $t \in [T]$ we have the difference of discrepancy bounded by

$$0 \leq \mathcal{L}_{\mathcal{D}_t}(f_t, \mathcal{P}_{J_t}(f_t)) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_t}(f_t, g) \leq \beta. \quad (\text{E.13})$$

Step 2: Bound the out-sample training error between updates.

Consider an update is executed at step $t+1$, it directly implies that $\mathcal{L}_{\mathcal{D}_t}(f_t, f_t) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_t}(f_t, g) > 4\beta$, while the latest update at step τ_t ensures that $\mathcal{L}_{\mathcal{D}_{\tau_t-1}}(f_{\tau_t}, f_{\tau_t}) - \inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}_{\tau_t-1}}(f_{\tau_t}, g) \leq \beta$, where τ_t is the pointer of the latest update. Combined the results above with (E.13), we have

$$\mathcal{L}_{\mathcal{D}_t}(f_t, f_t) - \mathcal{L}_{\mathcal{D}_t}(f_t, \mathcal{P}(f_t)) > 3\beta, \quad \mathcal{L}_{\mathcal{D}_{\tau_t-1}}(f_{\tau_t}, f_{\tau_t}) - \mathcal{L}_{\mathcal{D}_{\tau_t-1}}(f_{\tau_t}, \mathcal{P}(f_{\tau_t})) \leq \beta. \quad (\text{E.14})$$

It indicates that the sum of squared empirical discrepancy between two adjacent updates follows

$$\sum_{i=\tau_t}^t \|l_{f_t}(f_t, f_t, \zeta_t)\|_2^2 = \mathcal{L}_{\mathcal{D}_{\tau_t:t}}(f_t, f_t) - \mathcal{L}_{\mathcal{D}_{\tau_t:t}}(f_t, \mathcal{P}_{J_t}(f_t)) > 2\beta, \quad (\text{E.15})$$

where denote $\mathcal{D}_{\tau_t:t} = \mathcal{D}_t / \mathcal{D}_{\tau_t}$. Based on the similar concentration arguments as Lemma E.2, we have the out-sample training error between updates is bounded by $\sum_{i=\tau_t}^t \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_i, f_i, \zeta_i)]\|_2^2 > \beta$.

Step 3: Bound the switching cost under the transferability constraint.

Denote $b_1, \dots, b_{\mathcal{N}(T)}, b_{\mathcal{N}(T)+1}$ be the sequence of updated steps such that $\tau_t \in \{b_t\}$ for all $t \in [T]$, and we fix the recorder $b_1 = 1$ and $b_{\mathcal{N}(T)+1} = T + 1$. Note that based on (E.15), the sum of out-sample training error shall have a lower bound such that

$$\sum_{t=1}^T \|\mathbb{E}_{\zeta_t}[l_{f_t}(f_t, f_t, \zeta_t)]\|_2^2 = \sum_{u=1}^{\mathcal{N}(T)} \sum_{t=b_u}^{b_{u+1}-1} \|\mathbb{E}_{\zeta_t}[l_{f_t}(f_t, f_t, \zeta_t)]\|_2^2 \geq \mathcal{N}(T)\beta. \quad (\text{E.16})$$

Besides, note that the in-sample training error $\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_t)]\|_2^2 \leq \mathcal{O}(\beta)$ for all $t \in [T]$ and based on the definition of transferability coefficient κ_G (see Definition 3), we have

$$\sum_{t=1}^T \|\mathbb{E}_{\zeta_t}[l_{f_t}(f_t, f_t, \zeta_t)]\|_2^2 \leq \mathcal{O}(\kappa_G \max\{\beta, \text{sp}(V^*)^2\} \log T + T \log T \epsilon^2) \quad (\text{E.17})$$

Combine (E.16) and (E.17), it holds $\mathcal{N}(T) \leq \mathcal{O}(\kappa_G \log T + \beta^{-1} T \log T \epsilon^2)$ and finish the proof. \square

F Proof of Results about Complexity Measures

In this section, we provide the proof of results about the complexity metrics mentioned in Section 3.

F.1 Proof of Lemma 3.1

Proof of Lemma 3.1. Recall that the eluder dimension is defined over the function class following

$$\mathcal{X}_{\mathcal{H}} := \{X_{f,f'}(s, a) = (r_f + \mathbb{P}_{f'} V_f)(s, a) : f, f' \in \mathcal{H}\}.$$

Start with the transferability coefficient, the given condition can be written as

$$\begin{aligned} \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 &= \sum_{i=1}^{t-1} [(r_{f_t} + \mathbb{P}_{f_t} V_{f_t} - r_{f^*} + \mathbb{P}_{f^*} V_{f_t})(s_i, a_i)]^2 \\ &= \sum_{i=1}^{t-1} [(X_{f_t, f_t} - X_{f_t, f^*})(s_i, a_i)]^2 \leq \mathcal{O}(\beta), \end{aligned} \quad (\text{F.1})$$

for all $t \in [T]$. Denote $\mathcal{X}_{\mathcal{H}} - \mathcal{X}_{\mathcal{H}} = \{f - f' : f, f' \in \mathcal{X}_{\mathcal{H}}\}$, the generalized pigeon-hole principle (see Lemma H.2) suggests that if we take $\Gamma = \mathcal{D}_{\Delta}$ and $\phi_t = X_{f_t, f_t} - X_{f_t, f^*} \in \mathcal{X}_{\mathcal{H}} - \mathcal{X}_{\mathcal{H}}$, it holds

$$\begin{aligned} \sum_{i=1}^t \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_i, f_i, \zeta_i)]\|^2 &= \sum_{i=1}^t [(X_{f_i, f_i} - X_{f_i, f^*})(s_i, a_i)]^2 \\ &\leq \mathcal{O}(\dim_{\text{DE}}(\mathcal{X}_{\mathcal{H}} - \mathcal{X}_{\mathcal{H}}, \mathcal{D}_{\Delta}, \epsilon) \beta \log t + \text{sp}(V^*)^2 \min\{\dim_{\text{DE}}(\mathcal{X}_{\mathcal{H}} - \mathcal{X}_{\mathcal{H}}, \mathcal{D}_{\Delta}, \epsilon), t\} + t\epsilon^2) \\ &= \mathcal{O}(\dim_{\text{E}}(\mathcal{X}_{\mathcal{H}}, \epsilon) \beta \log t + \text{sp}(V^*)^2 \min\{\dim_{\text{E}}(\mathcal{X}_{\mathcal{H}}, \epsilon), t\} + t\epsilon^2), \end{aligned} \quad (\text{F.2})$$

where the last equation is based on the fact that $\dim_{\text{DE}}(\mathcal{X}_{\mathcal{H}} - \mathcal{X}_{\mathcal{H}}, \mathcal{D}_{\Delta}, \epsilon) = \dim_{\text{E}}(\mathcal{X}_{\mathcal{H}}, \epsilon)$ according to Definitions 5 and 7. Let $d_{\text{E}} = \dim_{\text{E}}(\mathcal{X}_{\mathcal{H}}, \epsilon)$, then we have $\kappa_{\text{G}} \leq d_{\text{E}}$. To show function class with a low Eluder dimension also shares a low dominance coefficient, use Lemma H.3 and it follows

$$\begin{aligned} \sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t)^2 &= \sum_{t=1}^T \|\mathbb{E}_{\zeta_t}[l_{f_t}(f_t, f_t, \zeta_t)]\|_2^2 = \sum_{t=1}^T [(X_{f_t, f_t} - X_{f_t, f^*})(s_t, a_t)]^2 \\ &\leq \left[d_{\text{DE}} (1 + \log T) \sum_{t=1}^T \sum_{i=1}^{t-1} [(X_{f_t, f_t} - X_{f_t, f^*})(s_i, a_i)]^2 \right]^{1/2} + \text{sp}(V^*) \min\{d_{\text{DE}}, T\} + T\epsilon \\ &\leq \left[d_{\text{E}} (1 + \log T) \sum_{t=1}^T \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \right]^{1/2} + \text{sp}(V^*) \min\{d_{\text{E}}, T\} + T\epsilon, \end{aligned}$$

by taking $\Gamma = \mathcal{D}_{\Delta}$, $\phi_t = X_{f_t, f_t} - X_{f_t, f^*}$, $\Phi = \mathcal{X}_{\mathcal{H}} - \mathcal{X}_{\mathcal{H}}$, $C \lesssim \text{sp}(V^*)$ and using the fact that $1 + \log T \leq 2 \log T$ and the equivalence of $\dim_{\text{DE}}(\mathcal{X}_{\mathcal{H}} - \mathcal{X}_{\mathcal{H}}, \mathcal{D}_{\Delta}, \epsilon)$ and $\dim_{\text{E}}(\mathcal{X}_{\mathcal{H}}, \epsilon)$. \square

F.2 Proof of Lemma 3.2

Proof of Lemma 3.2. Consider the Bellman discrepancy function, defined as

$$l_{f'}(f, g, \zeta_t) = Q_g(s_t, a_t) - r(s_t, a_t) - V_f(s_{t+1}) + J_g,$$

and the expectation is taken over the transition state s_{i+1} from $\mathbb{P}(\cdot|s_i, a_i)$ such that

$$\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 = \sum_{i=1}^{t-1} \mathcal{E}(f_i)(s_i, a_i)^2, \quad \forall t \in [T].$$

First, we're going to demonstrate the transferability. Note that the generalized pigeon-hole principle (see Lemma H.2) directly indicates that, given $\sum_{i=1}^{t-1} \|\mathcal{E}(f_t)(s_i, a_i)\|_2^2 \leq \mathcal{O}(\beta)$ holds true for all $t \in [T]$, if we take $\phi_t = \mathcal{E}(f_t)$, $\Phi = (I - \mathcal{T})\mathcal{H}$ and $\Gamma = \mathcal{D}_\Delta$, then we have

$$\sum_{i=1}^t \|\mathcal{E}(f_i)(s_i, a_i)\|_2^2 \leq \mathcal{O}(d_{\text{ABE}} \cdot \beta \log t + \text{sp}(V^*)^2 \min\{d_{\text{ABE}}, t\} + t\epsilon^2), \quad \forall t \in [T], \quad (\text{F.3})$$

and thus we upper bound the transferability coefficient by $\kappa_G \leq d_{\text{ABE}}$, where denote $d_{\text{ABE}} = \dim_{\text{ABE}}(\mathcal{H}, \epsilon)$. To demonstrate that the function class \mathcal{H} with a low ABE dimension also shares a low dominance coefficient, we can use Lemma H.3, which directly provides the result:

$$\sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t) \leq \left[2d_{\text{ABE}} \log T \sum_{t=1}^T \sum_{i=1}^{t-1} \|\mathcal{E}(f_t)(s_i, a_i)\|_2^2 \right]^{1/2} + \mathcal{O}(\text{sp}(V^*) \min\{d_{\text{ABE}}, T\}) + T\epsilon,$$

where we takes $\phi_t = \mathcal{E}(f_t)$, $\Phi = (I - \mathcal{T})\mathcal{H}$ and $\Gamma = \mathcal{D}_\Delta$ for Lemma H.3. Here, we also use the fact that $1 + \log T \leq 2 \log T$ and $\|\mathcal{E}(f_t)(s, a)\|_1 \lesssim \text{sp}(V^*)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. The inequality above demonstrates that the dominance coefficient is upper bounded by $d_G \leq 2d_{\text{ABE}} \cdot \log T$. \square

G Proof of Results for Concrete Examples

In this section, we provide detailed proofs of results for concrete examples in Appendix D.

G.1 Proof of Proposition D.1

Proof of Proposition D.1. To demonstrate that the linear FA has low AGECEC, we will prove that it can be characterized by a low ABE dimension, denoted as $\dim_{\text{ABE}}(\mathcal{H}, \epsilon)$, and the result can be directly drawn using Lemma 3.2. Recall definition of the ABE dimension, we assume that there exists constant $\epsilon' \geq \epsilon$, a sequence of distributions $\{\delta_{s_i, a_i}\}_{i \in [m]} \subseteq \mathcal{D}_\Delta$ and hypotheses $\{f_i\}_{i \in [m]} \subseteq \mathcal{H}$ such that

$$\sqrt{\sum_{i=1}^{t-1} [\mathcal{E}(f_t)(s_i, a_i)]^2} \leq \epsilon', \quad |\mathcal{E}(f_t)(s_t, a_t)| > \epsilon', \quad \forall t \in [m], \quad (\text{G.1})$$

based on Definitions 7-8. Next, we present a detailed discussion about the concrete problems, including linear AMDPs, AMDPs with linear Bellmen completion, and AMDPs with generalized linear

completion defined in Definition 10-12. We assume that the feature mapping is all d -dimensional, and we first illustrate that these problems share a similar Bellman error structure.

(i). **Linear AMPDs.** As defined in Definition 10, it implies that for any $f_t \in \mathcal{H}$

$$\begin{aligned}\mathcal{E}(f_t)(s_i, a_i) &= \phi(s_i, a_i)^\top \boldsymbol{\omega}_t - \phi(s_i, a_i)^\top \boldsymbol{\theta} - \phi(s_i, a_i)^\top \int_S V_t(s) d\boldsymbol{\mu}(s) + J_t \\ &= \left\langle \phi(s_i, a_i), \boldsymbol{\omega}_t - \boldsymbol{\theta} + \int_S V_t(s) d\boldsymbol{\mu}(s) + J_t \mathbf{e}_{(1)} \right\rangle,\end{aligned}\tag{G.2}$$

where denotes $\mathbf{e}_{(1)} = (1, 0, \dots, 0) \in \mathbb{R}^d$.

(ii). **AMPDs with linear Bellmen completion.** As a natural extension to the linear AMPDs, linear Bellmen completeness (see Definition 11) suggests that Bellman errors shall have the form:

$$\begin{aligned}\mathcal{E}(f_t)(s_i, a_i) &= \phi(s_i, a_i)^\top \boldsymbol{\omega}_t - \phi(s_i, a_i)^\top \mathcal{T}_{J_t}(\boldsymbol{\omega}_t) \\ &= \langle \phi(s_i, a_i), \boldsymbol{\omega}_t - \mathcal{T}_{J_t}(\boldsymbol{\omega}_t) \rangle.\end{aligned}\tag{G.3}$$

(iii). **AMPDs with generalized linear completion.** Moreover, AMPDs with generalized linear completion (see Definition 12) further extends the standard linear FA by introducing the link function, and the α -bi-Lipschitz continuity ensures the explorability of the problem. Detailedly,

$$\begin{aligned}\mathcal{E}(f_t)(s_i, a_i) &= \sigma(\phi(s_i, a_i)^\top \boldsymbol{\omega}_t) - \sigma(\phi(s_i, a_i)^\top \mathcal{T}_{J_t}(\boldsymbol{\omega}_t)) \\ &\in \left[\frac{1}{\alpha} \langle \phi(s_i, a_i), \boldsymbol{\omega}_t - \mathcal{T}_{J_t}(\boldsymbol{\omega}_t) \rangle, \alpha \langle \phi(s_i, a_i), \boldsymbol{\omega}_t - \mathcal{T}_{J_t}(\boldsymbol{\omega}_t) \rangle \right],\end{aligned}\tag{G.4}$$

where the bound is based on the definition of α -bi-Lipschitz continuity in (D.4). Based on the Lemma H.4 and arguments above, we're ready to provide a unified proof for both the linear FA and generalized linear FA. If we substitute the arguments (G.2), (G.3) and (G.4) back into (G.1), we have

$$\sqrt{\sum_{i=1}^{t-1} [\langle \phi(s_i, a_i), \boldsymbol{\omega}_t - \mathcal{T}_{J_t}(\boldsymbol{\omega}_t) \rangle]^2} \leq \alpha \epsilon', \quad |\langle \phi(s_t, a_t), \boldsymbol{\omega}_t - \mathcal{T}_{J_t}(\boldsymbol{\omega}_t) \rangle| > \frac{\epsilon'}{\alpha},\tag{G.5}$$

for all $t \in [T]$. Detailedly, we take $\alpha = 1$ for standard linear FA (e.g. linear AMPDs, AMPDs with linear Bellman completion), and denote α as the Lipschitz continuity coefficient for generalized linear FA (e.g. AMPDs with generalized linear Bellman completion). Based on the d -upper bound lemma (see Lemma H.4 for formal statement), if we take $\boldsymbol{\phi}_t = \phi(s_t, a_t)$, $\boldsymbol{\psi}_t = \boldsymbol{\omega}_t - \mathcal{T}_{J_t}(\boldsymbol{\omega}_t)$, $B_\phi = \sqrt{2}$, $B_\psi = \text{sp}(V^*)\sqrt{d}$, $\varepsilon = \epsilon$, $c_1 = \alpha$, $c_2 = \alpha^{-1}$, then the length of the sequence should be upper bounded by $m \leq \mathcal{O}(d \log(\text{sp}(V^*)\sqrt{d}/\epsilon))$ with a d -dimensional feature mapping. Thus, we have

$$\dim_{\text{ABE}}(\mathcal{H}, \epsilon) \leq \mathcal{O}\left(d \log(\text{sp}(V^*)\sqrt{d}/\epsilon)\right),$$

as the ABE dimension is defined as the longest sequence satisfying the condition in (G.5). Based on Lemma 3.2, $d_G \leq \mathcal{O}(d \log(\text{sp}(V^*)\sqrt{d}\epsilon^{-1}) \log T)$ and $\kappa_G \leq \mathcal{O}(d \log(\text{sp}(V^*)\sqrt{d}\epsilon^{-1}))$. \square

G.2 Proof of Proposition D.2

Proof of Proposition D.2. Consider the linear Q^*/V^* AMDPs, recall that the definition of linear Q^*/V^* indicates that the Bellman error can be written as below:

$$\begin{aligned}\mathcal{E}(f_t)(s_i, a_i) &= \phi(s_i, a_i)^\top \boldsymbol{\omega}_t - \mathbb{E}_{s_{i+1} \sim \mathbb{P}(\cdot | s_i, a_i)} [\boldsymbol{\psi}(s_{i+1})]^\top \boldsymbol{\theta}_t + J_t - r(s_i, a_i) \\ &= \phi(s_i, a_i)^\top \boldsymbol{\omega}_t - \mathbb{E}_{s_{i+1} \sim \mathbb{P}(\cdot | s_i, a_i)} [\boldsymbol{\psi}(s_{i+1})]^\top \boldsymbol{\theta}_t \\ &\quad - (Q^*(s_i, a_i) - \mathbb{E}_{s_{i+1} \sim \mathbb{P}(\cdot | s_i, a_i)} [V^*(s_{i+1})] - J^*) + J_t \quad (\text{see (2.1)}) \\ &= \left\langle \begin{bmatrix} \phi(s_i, a_i) \\ \mathbb{E}_{s_{i+1} \sim \mathbb{P}(\cdot | s_i, a_i)} [\boldsymbol{\psi}(s_{i+1})] \end{bmatrix}, \begin{bmatrix} \boldsymbol{\omega}_t - \boldsymbol{\omega}^* \\ \boldsymbol{\theta}^* - \boldsymbol{\theta}_t \end{bmatrix} + (J_t - J^*) \mathbf{e}_{(1)} \right\rangle, \quad (\text{G.6})\end{aligned}$$

where $\boldsymbol{\omega}^* \in \mathbb{R}^d$ and $\boldsymbol{\theta}^* \in \mathbb{R}^d$ denote the true optimal parameter, and $\mathbf{e}_{(1)} = (1, 0, \dots, 0) \in \mathbb{R}^d$. Following a similar argument in the proof of Proposition D.1, we can show that the linear Q^*/V^* AMDPs have a low ABE dimension with an equivalent $(d_1 + d_2)$ -dimensional compound feature mapping in (G.6). Based on Lemma 3.2, then we have it has low AGECE such that

$$d_G \leq \mathcal{O}(\underline{d} \log(\text{sp}(V^*) \sqrt{\underline{d}} \epsilon^{-1}) \log T), \quad \kappa_G \leq \mathcal{O}(\underline{d} \log(\text{sp}(V^*) \sqrt{\underline{d}} \epsilon^{-1}))$$

where denote $\underline{d} = d_1 + d_2$ be the sum of dimensions. \square

G.3 Proof of Proposition D.3

Proof of Proposition D.3. Similar to linear FA, we will prove that the kernel FA can be characterized by a low ABE dimension, denoted as $\dim_{\text{ABE}}(\mathcal{H}, \epsilon)$, and the result can be directly drawn from Lemma 3.2. Based on the definition of the ABE dimension, we assume that there exists constant $\epsilon' \geq \epsilon$, a sequence of distributions $\{\delta_{s_i, a_i}\}_{i \in [m]} \subseteq \mathcal{D}_\Delta$ and hypotheses $\{f_i\}_{i \in [m]} \subseteq \mathcal{H}$ such that

$$\sqrt{\sum_{i=1}^{t-1} [\mathcal{E}(f_t)(s_i, a_i)]^2} \leq \epsilon', \quad |\mathcal{E}(f_t)(s_t, a_t)| > \epsilon', \quad \forall t \in [m]. \quad (\text{G.7})$$

Suppose the kernel function class has a finite ϵ -effective dimension concerning the feature mapping ϕ . The existence of Bellman error $\mathcal{E}(f_t)$ is equivalent to the one of $W_t \in (\mathcal{W} - \mathcal{W})$:

$$\mathcal{E}(f_t)(\cdot, \cdot) = (Q_{f_t} - \mathcal{T}_{J_t} Q_{f_t})(\cdot, \cdot) = \langle \phi(\cdot, \cdot), \boldsymbol{\omega}_t - \boldsymbol{\omega}'_t \rangle_{\mathcal{K}} := \langle \phi(\cdot, \cdot), W_t \rangle_{\mathcal{K}}, \quad (\text{G.8})$$

where the second equation is based on the self-completeness assumption with kernel FA such that $\mathcal{G} = \mathcal{H}$. Denote $X_t = \phi(s_t, a_t)$, we can rewrite the condition in (G.7) as

$$\sqrt{\sum_{i=1}^{t-1} (X_i^\top W_t)^2} \leq \epsilon', \quad |X_t^\top W_t| > \epsilon', \quad \forall t \in [m]. \quad (\text{G.9})$$

Define $\Sigma_t = \sum_{i=1}^{t-1} X_i X_i^\top + \{\epsilon'^2/4R^2 \cdot \text{sp}(V^*)^2\} \cdot \mathbf{I}$, then $\|W_t\|_{\Sigma_t} \leq \sqrt{2}\epsilon'$ and $\epsilon' \leq \sqrt{2}\epsilon' \cdot \|X_t\|_{\Sigma_t^{-1}}$ using the Cauchy-Swartz inequality based on (G.9). Thus, we have $\|X_t\|_{\Sigma_t^{-1}}^2 \geq 0.5$ and we have

$$\begin{aligned} \sum_{t=1}^m \log \left(1 + \|X_t\|_{\Sigma_t^{-1}}^2 \right) &= \log \left(\det[\Sigma_{m+1}] / \det[\Sigma_1] \right) \\ &= \log \det \left[\mathbf{I} + \frac{4R^2 \cdot \text{sp}(V^*)^2}{\epsilon'^2} \sum_{t=1}^m X_t X_t^\top \right], \end{aligned} \quad (\text{G.10})$$

based on the matrix determinant lemma. Therefore, the (G.9) directly implies that

$$\frac{1}{e} \leq \log \frac{3}{2} \leq \frac{1}{m} \log \det \left[\mathbf{I} + \frac{4R^2 \text{sp}(V^*)^2}{\epsilon'^2} \sum_{t=1}^m X_t X_t^\top \right] \Rightarrow m \leq \dim_{\text{eff}}(\mathcal{X}, \epsilon/2\text{sp}(V^*)R).$$

Recall that the ϵ -effective dimension is defined as the minimum positive integer satisfying the condition. As the ABE dimension is defined as the length of the longest sequence satisfying (G.9), then we can bound ABE dimension by $\dim_{\text{ABE}}(\mathcal{H}, \epsilon) \leq \dim_{\text{eff}}(\mathcal{X}, \epsilon/2\text{sp}(V^*)R)$. Based on the Lemma 3.2, we have $d_G \leq \dim_{\text{eff}}(\mathcal{X}, \epsilon/2\text{sp}(V^*)R) \log T$ and $\kappa_G \leq \dim_{\text{eff}}(\mathcal{X}, \epsilon/2\text{sp}(V^*)R)$. \square

G.4 Proof of Proposition D.4

Proof of Proposition D.4. Note that expected discrepancy function follows: for any $t \in [T]$

$$\begin{aligned} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2 &= \boldsymbol{\theta}_t^\top \left(\boldsymbol{\psi}(s_i, a_i) + \int_{\mathcal{S}} \phi(s_i, a_i, s') V_{f_i}(s') ds' \right) - r(s_i, a_i) - \mathbb{E}_{\zeta_i}[V_{f_i}(s_{i+1})] \\ &= (\boldsymbol{\theta}_t - \boldsymbol{\theta}^*)^\top \left(\boldsymbol{\psi}(s_i, a_i) + \int_{\mathcal{S}} \phi(s_i, a_i, s') V_{f_i}(s') ds' \right). \end{aligned} \quad (\text{realizability})$$

Denote $\boldsymbol{\omega}_t = \boldsymbol{\theta}_t - \boldsymbol{\theta}^*$ and $X_t = \boldsymbol{\psi}(s_i, a_i) + \int_{\mathcal{S}} \phi(s_i, a_i, s') V_{f_i}(s') ds'$, we define that $\Sigma_t = \epsilon \mathbf{I} + \sum_{i=1}^{t-1} X_i X_i^\top$ given any constant $\epsilon > 0$. Note that

$$\begin{aligned} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_{\Sigma_t} &= \left[\epsilon \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2^2 + \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \right]^{1/2} \\ &\leq 2\sqrt{\epsilon} + \left[\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \right]^{1/2} \end{aligned} \quad (\text{G.11})$$

where we use $\|\boldsymbol{\theta}_t\| \leq 1$. Based on the elliptical potential lemma (see Lemma H.6), we have

$$\begin{aligned} \sum_{t=1}^T \|X_t\|_{\Sigma_t^{-1}} \wedge 1 &\leq \sum_{t=1}^T 2d \cdot \log \left(1 + d^{-1} \sum_{t=1}^T \|X_t\|_2 \right) \\ &\leq 2d \cdot \log \left(1 + (1 + \text{sp}(V^*)/2) \cdot T \left(\sqrt{d}\epsilon \right)^{-1} \right) := d(\epsilon), \end{aligned} \quad (\text{G.12})$$

where the last inequality results from scaling conditions $\|\phi\|_\infty \leq \sqrt{d}$, $\|\psi\|_\infty \leq \sqrt{d}$ and $\|V_f\|_\infty \leq \frac{1}{2}\text{sp}(V^*)$. Combine (G.11) and the fact that $\mathbf{1}(\|X_t\|_{\Sigma_t^{-1}} \geq 1) \leq \|X_t\|_{\Sigma_t^{-1}} \wedge 1$, then it holds

$$\sum_{t=1}^T \mathbf{1}(\|X_t\|_{\Sigma_t^{-1}} \geq 1) \leq \sum_{t=1}^T \|X_t\|_{\Sigma_t^{-1}} \wedge 1 \leq d(\epsilon). \quad (\text{G.13})$$

Dominance Note the sum of Bellman errors follows that

$$\begin{aligned} \sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t) &= \sum_{t=1}^T ((r_{f_t} + \mathbb{P}_{f_t} V_{f_t})(s_t, a_t) - (r_{f^*} + \mathbb{P}_{f^*} V_{f_t})(s_t, a_t)) \\ &= \sum_{t=1}^T (\theta_t - \theta^*)^\top \left(\psi(s_t, a_t) + \int_{\mathcal{S}} \phi(s_t, a_t, s') V_{f_t}(s') ds' \right) \\ &= \sum_{t=1}^T \omega_t^\top X_t \cdot \left(\mathbf{1}(\|X_t\|_{\Sigma_t^{-1}} \leq 1) + \mathbf{1}(\|X_t\|_{\Sigma_t^{-1}} > 1) \right) \\ &\leq \sum_{t=1}^T \omega_t^\top X_t \cdot \mathbf{1}(\|X_t\|_{\Sigma_t^{-1}} \leq 1) + (\text{sp}(V^*) + 2) \cdot \min\{d(\epsilon), T\} \\ &\leq \sum_{t=1}^T \|\omega_t\|_{\Sigma_t} \cdot (\|X_t\|_{\Sigma_t^{-1}} \wedge 1) + (\text{sp}(V^*) + 2) \cdot \min\{d(\epsilon), T\}, \end{aligned} \quad (\text{G.14})$$

where the first inequality uses (G.13) and the last inequality follows the Cauchy-Swartz inequality. Based on the inequality (G.11) and (G.12), we have

$$\begin{aligned} &\sum_{t=1}^T \|\omega_t\|_{\Sigma_t} \cdot (\|X_t\|_{\Sigma_t^{-1}} \wedge 1) \\ &\leq \sum_{t=1}^T \left(2\sqrt{\epsilon} + \left[\sum_{i=1}^{t-1} \|l_{f_i}(f_t, f_t, \zeta_i)\|_2^2 \right]^{1/2} \right) \cdot (\|X_t\|_{\Sigma_t^{-1}} \wedge 1) \\ &\leq \left[\sum_{t=1}^T 4\epsilon \right]^{1/2} \left[\sum_{t=1}^T \|X_t\|_{\Sigma_t^{-1}} \wedge 1 \right]^{1/2} + \left[d(\epsilon) \sum_{t=1}^T \sum_{i=1}^{t-1} \|l_{f_i}(f_t, f_t, \zeta_i)\|_2^2 \right]^{1/2} \\ &\leq 2\sqrt{T\epsilon \cdot \min\{d(\epsilon), T\}} + \left[d(\epsilon) \sum_{t=1}^T \sum_{i=1}^{t-1} \|l_{f_i}(f_t, f_t, \zeta_i)\|_2^2 \right]^{1/2}. \end{aligned} \quad (\text{G.15})$$

Plugging the result back into the inequality above, we conclude that

$$\begin{aligned}
\sum_{t=1}^T \mathcal{E}(f_t)(s_t, a_t) &\leq \left[d(\epsilon) \sum_{t=1}^T \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \right]^{1/2} \\
&\quad + 2\sqrt{T\epsilon \cdot \min\{d(\epsilon), T\}} + (\text{sp}(V^*) + 2) \min\{d(\epsilon), T\} \\
&\leq \left[d(\epsilon) \sum_{t=1}^T \sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \right]^{1/2} + (\text{sp}(V^*) + 3) \min\{d(\epsilon), T\} + T\epsilon, \quad (\text{G.16})
\end{aligned}$$

where the last inequality follows AM-GM inequality. Then, $d_G \leq \mathcal{O}(d \log(\text{sp}(V^*)T/\sqrt{d\epsilon}))$.

Transferability Given $\sum_{i=1}^{t-1} \|\mathbb{E}[l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 \leq \beta$ for all $t \in [T]$ and following the similar arguments in the proof of dominance, we have

$$\begin{aligned}
\sum_{t=1}^T \|\mathbb{E}[l_{f_t}(f_t, f_t, \zeta_t)]\|_2^2 &= \sum_{t=1}^T \left[(\boldsymbol{\theta}_t - \boldsymbol{\theta}^*)^\top \left(\boldsymbol{\psi}(s_t, a_t) + \int_S \boldsymbol{\phi}(s_t, a_t, s') V_{f_t}(s') ds' \right) \right]^2 \\
&\leq \sum_{t=1}^T (\boldsymbol{\omega}_t^\top X_t)^2 \cdot \mathbb{1} \left(\|X_t\|_{\Sigma_t^{-1}}^2 \leq 1 \right) + (\text{sp}(V^*) + 2)^2 \min\{d(\epsilon), T\} \\
&\leq \sum_{t=1}^T (\beta + 4\epsilon) \cdot \left(\|X_t\|_{\Sigma_t^{-1}}^2 \wedge 1 \right) + (\text{sp}(V^*) + 2)^2 \min\{d(\epsilon), T\} \\
&\leq d(\epsilon) \cdot \beta + (\text{sp}(V^*)^2 + 4 \cdot \text{sp}(V^*) + 6) \min\{d(\epsilon), T\} + 2T\epsilon^2, \quad (\text{G.17})
\end{aligned}$$

where we use a variant of (G.12) and (G.13), following that

$$\sum_{t=1}^T \mathbb{1} \left(\|X_t\|_{\Sigma_t^{-1}}^2 \geq 1 \right) \leq \sum_{t=1}^T \|X_t\|_{\Sigma_t^{-1}}^2 \wedge 1 \leq \sum_{t=1}^T \|X_t\|_{\Sigma_t^{-1}} \wedge 1 \leq d(\epsilon),$$

and the last inequality follows a similar proof as (G.15) and (G.16) based on Cauchy-Swartz and AM-GM inequality. Recall the definition of $d(\epsilon)$ in (G.12), then $\kappa_G \leq \mathcal{O}(d \log(\text{sp}(V^*)T/\sqrt{d\epsilon}))$. \square

G.5 Discussion about Performance on Concrete Examples

In this subsection, we demonstrate the detailed performance of LOOP under specific scenarios. Note that the algorithm achieves an $\tilde{\mathcal{O}}(\sqrt{T})$ regret, which is nearly minimax optimal, in both linear AMDP and linear mixture AMDP. We remark that existing algorithms were tailored for specific problems individually (Wei et al., 2021; Wu et al., 2022), LOOP offers a unified approach that covers them with comparable performance with an additional cost of generality for trade-off.

Linear AMDP Recall that AMDP falls into the linear function class, defined as $\mathcal{H}_{\text{Lin}} = \{ (Q, J) : Q(\cdot, \cdot) = \boldsymbol{\omega}^\top \boldsymbol{\phi}(\cdot, \cdot) \mid \|\boldsymbol{\omega}\|_2 \leq \frac{1}{2} \text{sp}(V^*) \sqrt{d}, J \in \mathbb{J}(\mathcal{H}) \}$. Consider the ρ -covering number, note that

$$|Q(s, a) - Q'(s, a)| \leq |(\boldsymbol{\omega} - \boldsymbol{\omega}')^\top \boldsymbol{\phi}(s, a)| \leq \sqrt{2} \|\boldsymbol{\omega} - \boldsymbol{\omega}'\|_1.$$

Based on the Lemma H.9, combine the fact that $|\mathbb{J}(\mathcal{H})| \leq 2$ and its ρ -covering number $\mathcal{N}_\rho(\mathbb{J}(\mathcal{H}))$ is at most $2\rho^{-1}$, we can get the log covering number of the hypotheses class \mathcal{H}_{Lin} is upper bounded by

$$\log \mathcal{N}_{\mathcal{H}}(\rho) \leq d \log \left(\text{sp}(V^*) 2^{\frac{3}{2}} d^{\frac{3}{2}} \rho^{-2} \right), \quad (\text{G.18})$$

by taking $\alpha = w$, $P = d$, $B = \text{sp}(V^*)\sqrt{d}/2$. Recall that the AGECEC for linear function approximation is bounded in Proposition D.1, and it holds that

$$d_G \leq \mathcal{O}(d \log(\text{sp}(V^*)\sqrt{d}\rho^{-1}) \log T), \quad \kappa_G \leq \mathcal{O}(d \log(\text{sp}(V^*)\sqrt{d}\rho^{-1})). \quad (\text{G.19})$$

Combine (G.18), (G.19) and the regret guarantee shown in Theorem 4.1, we get

$$\text{Reg}(T) \leq \mathcal{O}\left(\text{sp}(V^*) \max\{d_G, \kappa_G\} \sqrt{T \log(T \mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(1/T)/\delta) \text{sp}(V^*)}\right) \leq \tilde{\mathcal{O}}\left(\text{sp}(V^*)^{\frac{3}{2}} d^{\frac{3}{2}} \sqrt{T}\right).$$

For linear AMDPs, our method achieves a regret bound of $\tilde{\mathcal{O}}(\text{sp}(V^*)^{\frac{3}{2}} d^{\frac{3}{2}} \sqrt{T})$ for both linear and generalized linear AMDPs. In comparison, the FOPO algorithm (Wei et al., 2021) achieves the best-known regret bound of $\tilde{\mathcal{O}}(\text{sp}(V^*) d^{\frac{3}{2}} \sqrt{T})$ for linear AMDPs. However, our method incurs an additional $\text{sp}(V^*)^{\frac{1}{2}}$ in the regret bound, which is inevitable and is a trade-off for generalization.

Linear mixture Recall that the Proposition D.4 posits that AGECEC of the linear mixture problem satisfies that $\max\{\sqrt{d_G}, \kappa_G\} \leq \mathcal{O}(d\sqrt{\log T})$. Note that the hypotheses class is defined as

$$\mathcal{H}_{\text{LM}} = \{(\mathbb{P}, r) : \mathbb{P}(\cdot|s, a) = \boldsymbol{\theta}^\top \boldsymbol{\phi}(s, a, \cdot), \quad r(s, a) = \boldsymbol{\theta}^\top \boldsymbol{\psi}(s, a) \mid \|\boldsymbol{\theta}\|_2 \leq 1\}.$$

Consider the covering number note that both transition function and reward can be written as

$$\begin{aligned} \text{(i). } & |(\mathbb{P} - \mathbb{P}')(s'|s, a)| = |(\boldsymbol{\theta} - \boldsymbol{\theta}')^\top \boldsymbol{\phi}(s, a, s')| \leq \sqrt{d} \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_1, \\ \text{(ii). } & |(r - r')(s, a)| = |(\boldsymbol{\theta} - \boldsymbol{\theta}')^\top \boldsymbol{\psi}(s, a)| \leq \sqrt{d} \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_1. \end{aligned}$$

Based on the Lemma H.9, the log covering number of \mathcal{H}_{LM} is upper bounded by

$$\log \mathcal{N}_{\mathcal{H}}(\rho) \leq 2d \log \left(d^{\frac{3}{2}} \rho^{-1} \right),$$

by taking $\boldsymbol{\alpha} = \boldsymbol{\theta}$, $P = d$, $B = 1$. Combine results above and Theorem 4.1, we get

$$\text{Reg}(T) \leq \mathcal{O}\left(\text{sp}(V^*) \max\{d_G, \kappa_G\} \sqrt{T \log(T \mathcal{N}_{\mathcal{H} \cup \mathcal{G}}^2(1/T)/\delta) \text{sp}(V^*)}\right) \leq \tilde{\mathcal{O}}\left(\text{sp}(V^*)^{\frac{3}{2}} d^{\frac{3}{2}} \sqrt{T}\right).$$

which shares the same rate of cumulative regret as linear AMDPs. At our best knowledge, the UCRL2-VTR (Wu et al., 2022) achieves the best $\tilde{\mathcal{O}}(Dd\sqrt{T})$ regret for linear mixture AMDP, where D denotes the diameter under the communicating AMDP assumption and $\text{sp}(V^*) \leq D$ (Wang et al., 2022). In comparison, our method achieves $\tilde{\mathcal{O}}(\text{sp}(V^*)^{\frac{3}{2}} d^{\frac{3}{2}} \sqrt{T})$. The two algorithms are incomparable under different assumptions and both achieve a near minimax optimal regret at $\tilde{\mathcal{O}}(\sqrt{T})$.

H Technical Lemmas

In this section, we provide useful technical lemmas used in later theoretical analysis. Most are directly borrowed from existing works and proof of modified lemmas is provided in Section H.1.

Lemma H.1. Given function class Φ defined on \mathcal{X} , and a family of probability measures Γ over \mathcal{X} . Suppose sequence $\{\phi_k\}_{k=1}^K \subset \Phi$ and $\{\mu_k\}_{k=1}^K \subset \Gamma$ satisfy that for all $k \in [K]$, $\sum_{t=1}^{k-1} (\mathbb{E}_{\mu_t}[\phi_k])^2 \leq \beta$. Then, for all $k \in [K]$, we have

$$\sum_{t=1}^k \mathbb{1}(|\mathbb{E}_{\mu_t}[\phi_t]| > \epsilon) \leq \left(\frac{\beta}{\epsilon^2} + 1\right) \dim_{\text{DE}}(\Phi, \Pi, \epsilon).$$

Proof. See Lemma 43 of Jin et al. (2021) for detailed proof.

Lemma H.2 (Pigeon-hole principle). Given function class Φ defined on \mathcal{X} with $|\phi(x)| \leq C$ for all $\phi \in \Phi$ and $x \in \mathcal{X}$, and a family of probability measure over \mathcal{X} . Suppose sequence $\{\phi_k\}_{k=1}^K \subset \Phi$ and $\{\mu_k\}_{k=1}^K \subset \Gamma$ satisfy that for all $k \in [K]$, it holds $\sum_{t=1}^{k-1} (\mathbb{E}_{\mu_t}[\phi_k])^2 \leq \beta$. Let $d_{\text{DE}} = \dim_{\text{DE}}(\Phi, \Gamma, \epsilon)$ be the DE dimension, then for all $k \in [K]$ and $\epsilon > 0$, we have

$$\sum_{t=1}^k |\mathbb{E}_{\mu_t}[\phi_t]| \leq \mathcal{O}\left(\sqrt{d_{\text{DE}}\beta k} + \min\{k, d\}C + k\epsilon\right),$$

and

$$\sum_{t=1}^k \left[\mathbb{E}_{\mu_t}[\phi_t]\right]^2 \leq \mathcal{O}\left(d_{\text{DE}}\beta \log k + \min\{k, d\}C^2 + k\epsilon^2\right).$$

Proof. See Section H.1.1.

Lemma H.3. Given function class Φ defined on \mathcal{X} with $|\phi(x)| \leq C$ for all $\phi \in \Phi$ and $x \in \mathcal{X}$, and a family of probability measure over \mathcal{X} . Let $d_{\text{DE}} = \dim_{\text{DE}}(\Phi, \Gamma, \epsilon)$ be the DE dimension, then for all $k \in [K]$ and $\epsilon > 0$, we have

$$\sum_{t=1}^k |\mathbb{E}_{\mu_k}[\phi_k]| \leq \left[d_{\text{DE}} (1 + \log K) \sum_{k=1}^K \sum_{t=1}^{k-1} (\mathbb{E}_{\mu_t}[\phi_k])^2 \right]^{1/2} + \min\{d_{\text{DE}}, k\}C + k\epsilon.$$

Proof. See Section H.1.2.

Lemma H.4 (d -upper bound). Let Φ and Ψ be sets of d -dimensional vectors and $\|\phi\|_2 \leq B_\phi$, $\|\psi\|_2 \leq B_\psi$ for any $\phi \in \Phi$ and $\psi \in \Psi$. If there exists set (ϕ_1, \dots, ϕ_m) and (ψ_1, \dots, ψ_m) such that for all $t \in [m]$, $\sqrt{\sum_{k=1}^{t-1} \langle \phi_t, \psi_k \rangle^2} \leq c_1 \epsilon$ and $|\langle \phi_t, \psi_t \rangle| > c_2 \epsilon$, where $c_1 \geq c_2 > 0$ is a constant and $\epsilon > 0$, then the number of elements in set is bounded by $m \leq \mathcal{O}(d \log(B_\phi B_\psi / \epsilon))$.

Proof. See Section H.1.3.

Lemma H.5. For any sequence of positive reals x_1, \dots, x_m , it holds that

$$\frac{\sum_{i=1}^m x_i}{\sqrt{\sum_{i=1}^m i x_i^2}} \leq \sqrt{1 + \log n}.$$

Proof. See Lemma 6 in [Dann et al. \(2021\)](#) for detailed proof.

Lemma H.6. Let $\{x_i\}_{i \in [t]}$ be a sequence of vectors defined over Hilbert space \mathcal{X} . Let Λ_0 be a positive definite matrix and $\Lambda_t = \Lambda_0 + \sum_{i=1}^{t-1} x_i x_i^\top$. It holds that

$$\sum_{i=1}^t \|x_i\|_{\Lambda_t^{-1}}^2 \wedge 1 \leq 2 \log \left(\frac{\det(\Lambda_{t+1})}{\det(\Lambda_0)} \right).$$

Proof. See Elliptical Potential Lemma (EPL) in [Dani et al. \(2008\)](#) for a detailed proof.

Lemma H.7 (Freedman's inequality). Let X_1, \dots, X_T be a real-valued martingale difference sequence adapted to filtration $\{\mathcal{F}_t\}_{t=1}^T$. Assume for all $t \in [T]$ $X_t \leq R$, then for any $\eta \in (0, 1/R)$, with probability greater than $1 - \delta$

$$\sum_{t=1}^T X_t \leq \mathcal{O} \left(\eta \sum_{t=1}^T \mathbb{E}[X_t^2 | \mathcal{F}_t] + \frac{\log(1/\delta)}{\eta} \right),$$

Proof. See Lemma 7 in [Agarwal et al. \(2014\)](#) for detailed proof.

Lemma H.8 (Scaling lemma). Let $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ be a d -dimensional feature mapping, there exists an invertible linear transformation $A \in \mathbb{R}^{d \times d}$ such that for any bounded function $f : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ and $\mathbf{z} \in \mathbb{R}^d$ defined by

$$f(s, a) = \phi(s, a)^\top \mathbf{z},$$

we have $\|A\phi(s, a)\| \leq 1$ and $\|A^{-1}\mathbf{z}\| \leq \sup_{s,a} |f| \sqrt{d}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Proof. See Lemma 8 in [Wei et al. \(2021\)](#) for detailed proof.

In Theorem 4.1, the proved regret contains the logarithmic term of the $1/T$ -covering number of the function classes $\mathcal{N}_{\mathcal{H}}(1/T)$, which can be regarded as a surrogate cardinality of the function class \mathcal{H} . Here, we provide a formal definition of ρ -covering and the upper bound of ρ -covering number.

Definition 16 (ρ -covering). The ρ -covering number of a function class \mathcal{F} is the minimum integer t satisfying that there exists subset $\mathcal{F}' \subseteq \mathcal{F}$ with $|\mathcal{F}'| = t$ such that for any $f \in \mathcal{F}$ we can find a correspondence $f' \in \mathcal{F}'$ that it holds $\|f - f'\|_\infty \leq \rho$.

Lemma H.9 (ρ -covering number). Let \mathcal{F} be a function defined over \mathcal{X} that can be parametrized by $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_P) \in \mathbb{R}^P$ with $|\alpha_i| \leq B$ for all $i \in [P]$. Suppose that for any $f, f' \in \mathcal{F}$ it holds that $\sup_{x \in \mathcal{X}} |f(x) - f'(x)| \leq L \|\boldsymbol{\alpha} - \boldsymbol{\alpha}'\|_1$ and let $\mathcal{N}_{\mathcal{F}}(\rho)$ be the ρ -covering number of \mathcal{F} , then

$$\log \mathcal{N}_{\mathcal{F}}(\rho) \leq P \log \left(\frac{2BLP}{\rho} \right).$$

Proof. See Lemma 12 in [Wei et al. \(2021\)](#) for detailed proof.

H.1 Proof of Technical Lemmas

In this subsection, we present the proofs of technical auxiliary lemmas with modifications.

H.1.1 Proof of Lemma H.2

Proof of Lemma H.2. The first statement is directly from Lemma 41 in [Jin et al. \(2021\)](#), and the second statement follows a similar procedure as below. Note that Lemma H.1 suggests that

$$\sum_{t=1}^k \mathbb{1}\left([\mathbb{E}_{\mu_t}[\phi_t]]^2 > \epsilon^2\right) \leq \left(\frac{\beta}{\epsilon^2} + 1\right) \dim_{\text{DE}}(\Phi, \Gamma, \epsilon),$$

and note that the sum of squared expectation can be decomposed as

$$\begin{aligned} \sum_{t=1}^k [\mathbb{E}_{\mu_t}[\phi_t]]^2 &= \sum_{t=1}^k [\mathbb{E}_{\mu_t}[\phi_t]]^2 \mathbb{1}\left([\mathbb{E}_{\mu_t}[\phi_t]]^2 > \epsilon^2\right) + \sum_{t=1}^k [\mathbb{E}_{\mu_t}[\phi_t]]^2 \mathbb{1}\left([\mathbb{E}_{\mu_t}[\phi_t]]^2 \leq \epsilon^2\right) \\ &\leq \sum_{t=1}^k [\mathbb{E}_{\mu_t}[\phi_t]]^2 \mathbb{1}\left([\mathbb{E}_{\mu_t}[\phi_t]]^2 > \epsilon^2\right) + k\epsilon^2. \end{aligned} \quad (\text{H.1})$$

Assume sequence $[\mathbb{E}_{\mu_1}[\phi_1]]^2, \dots, [\mathbb{E}_{\mu_k}[\phi_k]]^2$ are sorted in the decreasing order and consider $t \in [k]$ such that $[\mathbb{E}_{\mu_t}[\phi_t]]^2 > \epsilon^2$, there exists a constant $\alpha \in (\epsilon^2, [\mathbb{E}_{\mu_t}[\phi_t]]^2)$ satisfying

$$t \leq \sum_{i=1}^k \mathbb{1}\left([\mathbb{E}_{\mu_i}[\phi_i]]^2 > \alpha\right) \leq \left(\frac{\beta}{\alpha} + 1\right) \dim_{\text{DE}}(\Phi, \Gamma, \sqrt{\alpha}) \leq \left(\frac{\beta}{\alpha} + 1\right) \dim_{\text{DE}}(\Phi, \Gamma, \epsilon),$$

where the last inequality is based on the fact that the DE dimension is monotonically decreasing in terms of ϵ as proposed in [Jin et al. \(2021\)](#). Denote $d_{\text{DE}} = \dim_{\text{DE}}(\Phi, \Gamma, \epsilon)$ and the inequality above implies that $\alpha \leq d_{\text{DE}}\beta/t - d$. Thus, we have $[\mathbb{E}_{\mu_t}[\phi_t]]^2 \leq d_{\text{DE}}\beta/t - d$. Beside, based on the definition we also have $[\mathbb{E}_{\mu_t}[\phi_t]]^2 \leq C^2$ and thus $[\mathbb{E}_{\mu_t}[\phi_t]]^2 \leq \min\{d_{\text{DE}}\beta/t - d, C^2\}$, then

$$\begin{aligned} \sum_{t=1}^k [\mathbb{E}_{\mu_t}[\phi_t]]^2 \mathbb{1}\left([\mathbb{E}_{\mu_t}[\phi_t]]^2 > \epsilon^2\right) &\leq \min\{d_{\text{DE}}, k\} C^2 + \sum_{t=d+1}^k \left(\frac{d_{\text{DE}}\beta}{t - d_{\text{DE}}}\right) \\ &\leq \min\{d_{\text{DE}}, k\} C^2 + d_{\text{DE}} \cdot \beta \int_0^k \frac{1}{t} dt \\ &\leq \min\{d_{\text{DE}}, k\} C^2 + d_{\text{DE}} \cdot \beta \log k. \end{aligned} \quad (\text{H.2})$$

Combine (H.1) and (H.2), then finishes the proof. \square

H.1.2 Proof of Lemma H.3

We remark that the proof provided in this subsection follows the almost same procedure as Lemma 3.16 in [Zhong et al. \(2022\)](#) with adjustment, and we preserve it for comprehension.

Proof of Lemma H.3. Denote $d_{\text{DE}} = \dim_{\text{DE}}(\Phi, \Gamma, \epsilon)$, $\hat{\epsilon}_{t,k} = |\mathbb{E}_{\mu_t}[\phi_k]|$ and $\epsilon_{t,k} = \hat{\epsilon}_{t,k} \mathbf{1}(\hat{\epsilon}_{t,k} > \epsilon)$ for $t, k \in [K]$, $\mu_t \in \Gamma$ and $\phi_k \in \Phi$. The proof follows the procedure. Consider K empty buckets B_0, \dots, B_{K-1} as initialization, and we examine $\epsilon_{k,k}$ one by one for all $k \in [K]$ as below:

Case 1 If $\epsilon_{k,k} = 0$, i.e., $\hat{\epsilon}_{k,k} \leq \epsilon$, then discard it.

Case 2 If $\epsilon_{k,k} > 0$, i.e., $\hat{\epsilon}_{k,k} > \epsilon$, at bucket j we add k into B_j if $\sum_{t \leq k-1, t \in B_j} (\epsilon_{t,k})^2 \leq (\epsilon_{k,k})^2$, otherwise we continue with the next bucket B_{j+1} .

Denote by b_k the index of bucket that at step k the non-zero $\epsilon_{k,k}$ falls in, i.e. $k \in B_{b_k}$. Based on the rule above, it holds that

$$\sum_{k=1}^K \sum_{t=1}^{k-1} (\epsilon_{t,k})^2 \geq \sum_{k=1}^K \sum_{0 \leq j \leq b_k-1, b_k \geq 1} \sum_{t \leq k-1, t \in B_j} (\epsilon_{t,k})^2 \geq \sum_{k=1}^K b_k \cdot (\epsilon_{k,k})^2,$$

where the first inequality arises from $\{t \in B_j : t \leq k-1, 0 \leq j \leq b_k-1, b_k \geq 1\} \subseteq [k-1]$ due to the discarding of the b_k th bucket, and the second equality directly follows the allocation rule such that $\sum_{t \leq k-1, t \in B_j} (\epsilon_{t,k})^2 \geq (\epsilon_{k,k})^2$ for any $j \leq b_k-1$. Recall that based on the definition of distributional eluder (DE) dimension, it is suggested the size $|B_j|$ is no larger than d_{DE} . Then,

$$\begin{aligned} \sum_{k=1}^K b_k (\epsilon_{k,k})^2 &= \sum_{j=1}^{K-1} j \sum_{t \in B_j} (\epsilon_{t,t})^2 && \text{(re-summation)} \\ &\geq \sum_{j=1}^{K-1} \frac{j}{|B_j|} \left(\sum_{t \in B_j} \epsilon_{t,t} \right)^2 \geq \sum_{j=1}^{K-1} \frac{j}{d_{\text{DE}}} \left(\sum_{t \in B_j} \epsilon_{t,t} \right)^2 && (|B_j| \leq d_{\text{DE}}) \\ &\geq (d_{\text{DE}} (1 + \log K))^{-1} \left(\sum_{j=1}^{K-1} \sum_{t \in B_j} \epsilon_{t,t} \right)^2 = (d_{\text{DE}} (1 + \log K))^{-1} \left(\sum_{t \in [K] \setminus B_0} \epsilon_{t,t} \right)^2, && \text{(H.3)} \end{aligned}$$

where the second inequality follows Lemma H.5. Combine the (H.1.2) and (H.3) above, we have

$$\begin{aligned} \sum_{k=1}^K \hat{\epsilon}_{k,k} &\leq \sum_{k=1}^K \epsilon_{k,k} + K\epsilon \leq \sum_{t \in [K] \setminus B_0} \epsilon_{t,t} + \min\{d_{\text{DE}}, K\} \|\phi\|_{\infty} + K\epsilon \\ &\leq \left[d_{\text{DE}} (1 + \log K) \sum_{k=1}^K \sum_{t=1}^{k-1} (\epsilon_{t,k})^2 \right]^{1/2} + \min\{d_{\text{DE}}, K\} C + K\epsilon \\ &\leq \left[d_{\text{DE}} (1 + \log K) \sum_{k=1}^K \sum_{t=1}^{k-1} (\hat{\epsilon}_{t,k})^2 \right]^{1/2} + \min\{d_{\text{DE}}, K\} C + K\epsilon. \end{aligned}$$

Substitute the definition $\hat{\epsilon}_{t,k} = |\mathbb{E}_{\mu_t}[\phi_k]|$ back into the inequality, then finishes the proof. \square

H.1.3 Proof of Lemma H.4

Proof of Lemma H.4. For notation simplicity, denote $\Lambda_t = \sum_{k=1}^{t-1} \psi_k \psi_k^\top + \frac{\varepsilon^2}{B_\phi^2} \cdot \mathbf{I}$, then for all $t \in [m]$ we have $\|\phi_t\|_{\Lambda_t} \leq \sqrt{\sum_{k=1}^{t-1} (\phi_t^\top \psi_k)^2 + \frac{\varepsilon^2}{B_\phi^2} \|\phi_t\|_2^2} = \sqrt{c_1^2 + 1} \varepsilon$ based on the given condition. Using the Cauchy-Swartz inequality and results above, then it holds $\|\psi_t\|_{\Lambda_t^{-1}} \geq |\langle \phi_t, \psi_t \rangle| / \|\phi_t\|_{\Lambda_t} = c_2 / \sqrt{c_1^2 + 1}$. On one hand, the matrix determinant lemma ensures that

$$\det(\Lambda_m) = \det(\Lambda_0) \cdot \prod_{t=1}^{m-1} (1 + \|\psi_t\|_{\Lambda_t^{-1}}^2) \geq (1 + c_2^2 / (1 + c_1^2))^{m-1} (\varepsilon^2 / B_\phi^2)^d. \quad (\text{H.4})$$

On the other hand, according to the definition of Λ_t , we have

$$\det(\Lambda_m) \leq (\text{trace}(\Lambda_m) / d)^d \leq \left(\sum_{k=1}^{t-1} \|\psi_k\|_2^2 / d + \varepsilon^2 / B_\phi^2 \right)^d \leq (B_\psi^2 (m-1) / d + \varepsilon^2 / B_\phi^2)^d. \quad (\text{H.5})$$

Combine (H.4) and (H.5), if we take logarithms at both sides, then we have

$$m \leq 1 + d \log \left(\frac{B_\phi^2 B_\psi^2 (m-1)}{d \varepsilon^2} + 1 \right) / \log \left(1 + \frac{c_2^2}{1 + c_1^2} \right).$$

After simple calculations, we can obtain that m is upper bounded by $\mathcal{O}(d \log(B_\phi B_\psi / \varepsilon))$. \square

I Supplementary Discussions

I.1 Proof Sketch of MLE-based Results

In this section, we provide the proof sketch of Theorem C.1. First, we introduce several useful lemmas, which is the variant of ones in Appendix E for MLE-based problems, and most have been fully researched in Liu et al. (2022, 2023a); Xiong et al. (2023). As there's no significant technical gap between the episodic and average-reward for model-based problems, we only provide a sketch.

Lemma I.1 (Akin to Lemma E.1). Under Assumptions 1-2, MLE-LOOP is an optimistic algorithm such that it ensures $J_t \geq J^*$ for all $t \in [T]$ with probability greater than $1 - \delta$.

Proof Sketch of Lemma I.1. See Proposition 13 in Liu et al. (2022) with slight modifications. \square

Lemma I.2 (Akin to Lemma E.2). For fixed $\rho > 0$ and pre-determined optimistic parameter $\beta = c(\log(T\mathcal{B}_\mathcal{H}(\rho)/\delta) + T\rho)$ where constant $c > 0$, it holds that

$$\sum_{i=1}^{t-1} \|\mathbb{E}_{\zeta_i} [l_{f_i}(f_t, f_t, \zeta_i)]\|_2^2 = \sum_{i=1}^{t-1} \text{TV}(\mathbb{P}_{f_t}(\cdot | s_i, a_i), \mathbb{P}_{f^*}(\cdot | s_i, a_i))^2 \leq \mathcal{O}(\beta), \quad (\text{I.1})$$

for all $t \in [T]$ with probability greater than $1 - \delta$.

Proof Sketch of Lemma I.2. See Proposition 14 in Liu et al. (2022) with slight modifications. \square

Lemma I.3 (Akin to Lemma E.3). Let $\mathcal{N}(T)$ be the switching cost with time horizon T , given fixed covering coefficient $\rho > 0$ and pre-determined optimistic parameter $\beta = c(\log(T\mathcal{B}_{\mathcal{H}}(\rho)/\delta) + T\rho)$ where c is a large enough constant, with probability greater than $1 - 2\delta$ we have

$$\mathcal{N}(T) \leq \mathcal{O}(\kappa_G \cdot \text{poly}(\log T) + \beta^{-1}T \log T \epsilon^2),$$

where κ_G is the transferability coefficient with respect to $\text{MLE-AGEC}(\mathcal{H}, \{l_{f'}\}, \epsilon)$.

Proof Sketch of Lemma I.3. The proof is almost the same as Lemma E.3.

Step 1: Bound the difference of discrepancy between the minimizer and f^* .

As proposed in Proposition 14, Liu et al. (2022), with a high probability it holds that

$$0 \leq \sum_{i=1}^t \text{TV}(\mathbb{P}_{f^*}(\cdot|s_i, a_i), \mathbb{P}_{g_i}(\cdot|s_i, a_i)) \leq \sqrt{\beta t}, \quad \forall t \in [T]. \quad (\text{I.2})$$

Step 2: Bound the expected discrepancy between updates.

Note that for all $t + 1 \in [T]$, the update happens only if

$$\sum_{i=1}^t \text{TV}(\mathbb{P}_{f_t}(\cdot|s_i, a_i), \mathbb{P}_{g_i}(\cdot|s_i, a_i)) > 3\sqrt{\beta t}. \quad (\text{I.3})$$

Combine the (I.2) and (I.3) above, and apply the triangle inequality, we have

$$\begin{aligned} & \sum_{i=1}^t \text{TV}(\mathbb{P}_{f_t}(\cdot|s_i, a_i), \mathbb{P}_{f^*}(\cdot|s_i, a_i)) \\ & \geq \sum_{i=1}^t \text{TV}(\mathbb{P}_{f_t}(\cdot|s_i, a_i), \mathbb{P}_{g_t}(\cdot|s_i, a_i)) - \text{TV}(\mathbb{P}_{f^*}(\cdot|s_i, a_i), \mathbb{P}_{g_t}(\cdot|s_i, a_i)) \geq 2\sqrt{\beta t}. \end{aligned}$$

and the construction of confidence set ensures that $\sum_{i=1}^{T_t} \text{TV}(\mathbb{P}_{f_t}(\cdot|s_i, a_i), \mathbb{P}_{f^*}(\cdot|s_i, a_i)) \leq \sqrt{\beta t}$ with high probability (Liu et al., 2022, Proposition 14). Recall the definition of MLE-transferability coefficient, then the switching cost can be bounded following the same argument in Lemma E.3. \square

Proof Sketch of Theorem C.1. Recall that

$$\text{Reg}(T) \leq \underbrace{\sum_{i=1}^T \mathcal{E}(f_t)(s_t, a_t)}_{\text{Bellman error}} + \underbrace{\sum_{t=1}^T \left(\mathbb{E}_{s_{t+1} \sim \mathbb{P}(\cdot|s_t, a_t)}[V_t(s_{t+1})] - V_t(s_t) \right)}_{\text{Realization error}}, \quad (\text{I.4})$$

where the inequality follows the optimism in Lemma I.1. Combine Lemma I.2, Lemma I.3 and the definition of MLE-AGEC (see Definition 9), we can easily finish the proof of regret. \square

Algorithm 3 Extended Value Iteration (EVI)

Input: hypothesis $f = (\mathbb{P}_f, r_f)$, desired accuracy level ϵ .

Initialize: $V^{(0)}(s) = 0$ for all $s \in \mathcal{S}$, $J^{(0)} = 0$ and counter $i = 0$.

- 1: **repeat**
 - 2: **for** $s \in \mathcal{S}$ and $a \in \mathcal{A}$ **do**
 - 3: Set $Q^{(i)}(s, a) \leftarrow r_f(s, a) + \mathbb{E}_{s' \sim \mathbb{P}_f(s, a)}[V^{(i)}(s')] - J^{(i)}$
 - 4: Update $V^{(i+1)}(s) \leftarrow \max_{a \in \mathcal{A}} Q^{(i)}(s, a)$
 - 5: Update counter $i \leftarrow i + 1$
 - 6: **until** $\max_{s \in \mathcal{S}} \{V^{(i+1)}(s) - V^{(i)}(s)\} - \min_{s \in \mathcal{S}} \{V^{(i+1)}(s) - V^{(i)}(s)\} \leq \epsilon$
-

I.2 Extended Value Iteration (EVI) for Model-Based Hypotheses

In model-based problems, the discrepancy function sometimes relies on the optimal state bias function V_f and optimal average-reward J_f (see linear mixture model in Section D). In this section, we provide an algorithm, extended value iteration (EVI) proposed in Auer et al. (2008), to output the optimal function and average-reward under given a model-based hypothesis $f = (\mathbb{P}_f, r_f)$. See Algorithm 3 for complete pseudocode. The convergence of EVI is guaranteed by the theorem below.

Theorem I.4. Under Assumption 1, there exists a unique centralized solution pair (Q^*, J^*) to the Bellman optimality equation for any AMDP \mathcal{M}_f characterized by hypothesis $f \in \mathcal{H}$. Then, if the extended value iteration (EVI) is stopped under the condition that

$$\max_{s \in \mathcal{S}} \{V^{(i+1)}(s) - V^{(i)}(s)\} - \min_{s \in \mathcal{S}} \{V^{(i+1)}(s) - V^{(i)}(s)\} \leq \epsilon,$$

then the achieved greedy policy $\pi^{(i)}$ is ϵ -optimal such that $J_{\mathcal{M}_f}^{\pi^{(i)}} \geq J_{\mathcal{M}_f}^* + \epsilon$.

Proof Sketch: See Theorem 12 in Auer et al. (2008).