
Boosting for Reinforcement Learning in Structured MDPs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Boosting is a powerful machine learning technique that constructs a strong learner
2 by sequentially combining weak learners, each of which performs only slightly
3 better than random. Recent work has adapted boosting to reinforcement learning
4 and established global convergence guarantees under the assumption of access to
5 a multiplicative weak learner (Brukhim et al., 2022). These guarantees critically
6 depend on occupancy mismatch terms relative to the optimal policy, however
7 the mismatch ratio between the boosted policy class and the optimal policy can
8 become unbounded unless the policy class ensures sufficient state-space coverage.
9 In this work, we remove this assumption and show that, whenever weak learning
10 is feasible, boosting can achieve convergence guarantees that depend only on the
11 intrinsic complexity of the underlying Markov Decision Process.

12 1 Introduction

13 In reinforcement learning, an agent interacts with an environment, specified by a Markov Decision
14 Process (MDP), with the objective of learning a policy for interaction that maximizes its expected
15 cumulative reward. The computational and statistical complexity of this task depends heavily on
16 the structure of the underlying MDP, and while computationally efficient algorithms that provably
17 learn an optimal policy are known for MDPs with simple structure, such as tabular or low-rank
18 MDPs (Brafman and Tennenholtz, 2002; Kearns and Singh, 2002; Jin et al., 2020), such algorithms
19 are known not to exist for some more expressive environments like block MDPs (under reasonable
20 cryptographic assumptions (Kane et al., 2022)).

21 Nonetheless, heuristic methods like policy gradient methods and deep-Q learning have been used in
22 practice to obtain impressive performance in complex environments such as MDPs with continuous
23 state or action spaces. This observation motivates the study of ensembling approaches for reinforce-
24 ment learning – in cases where efficient, provable optimization methods are not known to exist, but
25 policies that improve over random actions can be efficiently learned, how can we effectively generate
26 and ensemble a collection of “rule of thumb” policies to obtain something close to optimal?

27 Prior work of Brukhim et al. (2022) proposed a boosting-style approach to ensembling for RL,
28 inspired by boosting techniques from supervised learning (Freund, 1995). Assuming the existence of
29 a weak supervised learner for a parameterized policy class – a learner that outputs a policy somewhat
30 better than taking random actions, but potentially far from optimal – they give an algorithm to
31 iteratively refine a policy via aggregation with new policies generated by the weak learner. To ensure
32 each update improves the current policy, the reward function of the underlying MDP for the RL
33 instance is modified at each iteration to encourage the weak learner to output policies that capture
34 residual reward missed by the current policy.

35 The convergence of their algorithm to an approximately optimal policy is *independent of the size of*
36 *the state space* of the underlying MDP, and so is applicable to environments far more complex than

37 simple tabular MDPs. However, their convergence guarantees require two assumptions. The first is a
38 reasonable assumption on the expressibility of an optimal policy as a function of policies from the
39 class parameterizing the weak learner. The second is a common, but strong assumption arising from
40 off-policy learning – that the worst case distribution mismatch between the state visitation distribution
41 of the optimal policy and any policy in the parameterized class is bounded.

42 In this work, we aim to remove this second assumption by providing an alternative approach to
43 boosting in RL that explicitly encourages exploration. By iteratively altering the reward function of
44 the underlying MDP to promote visitation of new states, our algorithm converges independent of
45 distribution mismatch measures, guaranteeing convergence even if the starting state distribution and
46 that of the optimal policy are disjoint.

47 1.1 Our Results

48 We present a boosting approach for reinforcement learning that removes the dependence on occupancy
49 mismatch to the optimal policy present in [Brukhim et al. \(2022\)](#). In order to achieve this independence,
50 our algorithms modify the reward function of the underlying MDP to add a bonus term that rewards
51 visitation of under-explored states, using an approach that builds on the L_∞ and pushforward
52 coverability objectives of [Amortila et al. \(2024\)](#). By shifting the dependency from occupancy
53 mismatch to the less restrictive notion of coverability, our framework yields boosting-based RL
54 guarantees with sample complexity that scales polynomially with the MDP’s intrinsic parameters.
55 This weakens the coverage assumptions on the policy class required by the prior RL boosting approach
56 in [Brukhim et al. \(2022\)](#).

57 Our algorithms construct a collection of policies by iteratively invoking the boosting framework of
58 [Brukhim et al. \(2022\)](#). To avoid introducing a dependence on distribution mismatch for the target
59 MDP and policy class, we only invoke the boosting subroutine with a modified policy class and MDP
60 which guarantees efficient Frank-Wolfe style best-iterate convergence (for the modified MDP and
61 policy class). The modified MDP mimics the target MDP for states with good coverage under the
62 current collection of policies, so to ensure convergence for the target MDP, it suffices to show that
63 the modified MDP converges to something that well-approximates the target. We guarantee this by
64 modifying the reward function within each call to the boosting subroutine with an exploration bonus
65 using coverability notions from [Amortila et al. \(2024\)](#), ensuring that reachable states of the MDP are
66 visited with sufficiently high probability.

67 1.2 Paper organization

68 We include preliminaries on RL in Section 2, including pseudocode for the boosting framework
69 of [Brukhim et al. \(2022\)](#) which our algorithms use as a subroutine, as well as relevant concepts.
70 In Section 3 we present our two main boosting algorithms, and recall the L_∞ (Section 3.1) and
71 pushforward (Section 3.2) coverability notions from [Amortila et al. \(2024\)](#) that we use to define our
72 exploration bonuses. We summarize our contributions and future directions in Section 4. A glossary
73 of notation can be found in Appendix A. We discuss related work not covered in the introduction in
74 Appendix B. We provide the setup for the extended MDP and truncated policy class in Appendix C.
75 The complete proofs of the main text are in Appendix D and Appendix E. We provide a discussion
76 and results of using boosting for the reward-free setting in Appendix F. An adaptation of the RL
77 boosting framework of [Brukhim et al. \(2022\)](#) to episodic MDPs is given in Appendix G. Finally, the
78 results of empirical validation of our theoretical results can be found in Appendix K.

79 2 Preliminaries

80 **Markov Decision Process.** We consider an infinite-horizon discounted Markov Decision Process
81 (MDP) defined by the tuple $\mathcal{M} = (\mathcal{X}, \mathcal{A}, P, r, \gamma, d_0)$, where \mathcal{S} and \mathcal{A} are the state and action spaces,
82 respectively. The transition kernel $P(x' | x, a)$ denotes the probability of transitioning to state x'
83 after taking action a in state x . The reward function is given by $r : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$, $\gamma \in [0, 1)$ is the
84 discount factor, and $d_0 \in \Delta(\mathcal{X})$ is the initial state distribution.

85 An agent’s behavior is characterized by a stochastic policy $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$. The interaction between
86 π and \mathcal{M} induces a distribution over trajectories $\tau = (x_0, a_0, x_1, a_1, \dots)$, defined by $x_0 \sim d_0$,
87 $a_t \sim \pi(\cdot | x_t)$, and $x_{t+1} \sim P(\cdot | x_t, a_t)$. We define the state-action value function $Q^\pi(x, a)$ and the

88 state value function $V^\pi(x)$ as:

$$Q^\pi(x, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r(x_t, a_t) \mid x_0 = x, a_0 = a \right], \quad V^\pi(x) = \mathbb{E}_{a \sim \pi(\cdot|x)} [Q^\pi(x, a)].$$

89 For any distribution d_0 over \mathcal{X} , we denote $V_{d_0}^\pi = \mathbb{E}_{x \sim d_0} [V^\pi(x)]$. When the distribution is omitted,
90 V^π refers to $V_{d_0}^\pi$, and we denote the optimal value as $V^* = \max_\pi V^\pi$.

91 The policy π also induces a normalized discounted state occupancy measure $d_{d_0}^\pi \in \Delta(\mathcal{X})$, defined as:

$$d_{d_0}^\pi(x) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(x_t = x \mid \pi, x_0 \sim d_0). \quad (1)$$

92 Equivalently, the value function can be expressed via this occupancy measure as $V_{d_0}^\pi =$
93 $\frac{1}{1-\gamma} \mathbb{E}_{x \sim d_{d_0}^\pi, a \sim \pi(\cdot|x)} [r(x, a)]$.

94 **Modes of Accessing the MDP.** We consider the *episodic rollout setting* in which the algorithm can
95 execute a policy, terminate the rollout at any time, and restart from the initial states, and can repeat
96 this process multiple times.

97 In this work, we consider large-scale MDPs where the state space \mathcal{X} can be very large, or even infinite.
98 Such settings are challenging and statistically intractable in general (Krishnamurthy et al., 2016).
99 However, in many applications, the environment admits some underlying structure that enables both
100 computationally and statistically efficient algorithms; examples of such environments include block
101 MDPs (Misra et al., 2020; Du et al., 2019) and linear/low-rank MDPs (Jiang et al., 2017; Agarwal
102 et al., 2020). We provide a brief description of these structured MDPs in the following paragraphs.

103 **Block MDPs.** A block MDP models an environment where a large state space \mathcal{X} can be summarized
104 into a finite latent state \mathcal{S} . Specifically, block MDP assumes there exists an emission function
105 $q : \mathcal{S} \rightarrow \Delta(\mathcal{X})$ that maps the latent states \mathcal{S} to a distribution over \mathcal{X} . The agent interacts with the
106 environment by repeatedly generating trajectories $(s_0, x_0, a_0, r_0, s_1, x_1, a_1, r_1, \dots)$, where $s_0 \sim d_0$,
107 $s_{t+1} \sim P(\cdot \mid s_t, a_t)$, $x_t \sim q(s_t)$. The agent can observe x_t , but the latent states s_t are hidden from
108 the agent. Tabular MDP is a special case of block MDP where the observation space coincides with
109 the state space, i.e., $\mathcal{X} = \mathcal{S}$.

110 **Low-rank MDPs.** A low-rank MDP with dimension d is an MDP where the transition kernel
111 admits a low-rank factorization $\phi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d$ and $\mu : \mathcal{X} \rightarrow \mathbb{R}^d$ such that $\forall x, x' \in \mathcal{X}, a \in$
112 $\mathcal{A} : P(x' \mid x, a) = \langle \phi(x, a), \mu(x') \rangle$. This structure enables efficient representation and learning in
113 environments with large or continuous state spaces.

114 2.1 Boosting Framework (Brukhim et al., 2022)

115 In this section, we restate the key definitions and the boosting procedure in Brukhim et al. (2022).

116 **Definition 2.1** (Distribution Mismatch). Let $\pi^* = \operatorname{argmax}_\pi V^\pi$ denote the optimal policy. We define
117 the following distribution mismatch coefficients: $C_\infty = \max_{\pi \in \Pi} \left\| d^{\pi^*} / d^\pi \right\|_\infty$.

118 A natural approach to defining a weak learner is to require performance that is marginally better
119 than that of a random policy. In the definition below, let π_{Rand} be a uniformly random policy, i.e.
120 $\forall(x, a) \in \mathcal{X} \times \mathcal{A}, \pi_{\text{Rand}}(a \mid x) = 1/|\mathcal{A}|$.

121 **Definition 2.2** (Weak Supervised Learner). Let $\alpha \in (0, 1]$. Consider a class \mathcal{L} of linear loss functions
122 $\ell : \mathbb{R}^{\mathcal{A}} \rightarrow \mathbb{R}$, a family \mathbb{D} of distributions that are supported over $\mathcal{S} \times \mathcal{L}$, and policy class Π . A weak
123 supervised learning algorithm, for every $\varepsilon_w, \delta > 0$, given $m(\varepsilon_w, \delta) = \frac{\log |\mathcal{W}|}{\varepsilon_w^2} \log \frac{1}{\delta}$ samples D_m
124 from any distribution $\mathcal{D} \in \mathbb{D}$ outputs a policy $\mathcal{W}(D_m) \in \Pi$ such that with probability $1 - \delta$,

$$\mathbb{E}_{(s, \ell) \sim \mathcal{D}} [\ell(\mathcal{W}(D_m))] \leq \alpha \min_{\pi^* \in \Pi} \mathbb{E}_{(s, \ell) \sim \mathcal{D}} [\ell(\pi^*(s))] + (1 - \alpha) \mathbb{E}_{(s, \ell) \sim \mathcal{D}} [\ell(\pi_{\text{Rand}}(s))] + \varepsilon_w.$$

125 **Definition 2.3** (Policy Projection). Given $\tilde{\pi} : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{A}|}$, define a projected policy $\pi = \Gamma[\tilde{\pi}]$ to be a
126 policy such that simultaneously for all $x \in \mathcal{X}$, it holds that $\pi(\cdot \mid x) = \Gamma[\tilde{\pi}(x)]$.

127 Due to the multiplicative factor α , the policy obtained from the second boosting procedure (Internal
 128 Boost; [Brukhim et al. \(2022, Algorithm 2\)](#)) might not be a convex combination of weak policies.
 129 The projection operator Γ ensures that the policy obtained from the boosting procedure is valid. The
 130 projection can be computed efficiently using a water-filling algorithm ([Duchi et al., 2008](#)), which
 131 requires sorting $|\mathcal{A}|$ numbers ([Brukhim et al., 2022, Claim 12](#)).

132 The policy obtained from the RL boosting algorithm of [Brukhim et al. \(2022\)](#) can be presented as a
 133 two-level tree, which we restate below.

134 **Definition 2.4** (Policy Tree). A Policy Tree $\mathbb{P} \subseteq \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$ with respect to $\Pi \subseteq \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$ some
 135 base policy class, and $N, T \in \mathbb{N}$, is a linear combination of T projected policies $\Gamma[\tilde{\pi}]$, where each $\tilde{\pi}$
 136 is a linear combination of N base policies $\pi \in \Pi$.

137 **Definition 2.5** (Policy Completeness). For any initial state distribution μ , and policy classes Π, \mathbb{P} ,
 138 define

$$\mathcal{E}_{\mu}(\mathbb{P}, \Pi) = \max_{\pi \in \mathbb{P}} \min_{\pi^* \in \Pi} \mathbb{E}_{x \sim d_{\mu}^{\pi}} \left[\max_{a \in \mathcal{A}} Q^{\pi}(x, a) - Q^{\pi}(x, \cdot)^{\top} \pi^*(\cdot | x) \right].$$

139 Policy completeness quantifies how well the restricted policy class Π can approximate the greedy
 140 policy induced by the action-value function of any policy in the larger class \mathbb{P} , as measured under the
 141 state distribution d_{μ}^{π} .

We highlight the core boosting framework of [Brukhim et al. \(2022\)](#) in [Algorithm 1](#), and refer to the
[Appendix L.1](#) for their remaining subroutines and discussion. At a high level, the framework employs
 a two-level boosting scheme. The outer level follows a CPI procedure ([Kakade and Langford, 2002](#)),
 which aggregates policies produced by an inner boosting routine (Internal Boost). The Internal Boost
 is designed to return a policy π'_t ([Line 4](#) of [Algorithm 1](#)) that approximately closes the best policy
 under the state distribution $d^{\pi_{t-1}}$. Informally,

$$\mathbb{E}_{x \sim d^{\pi_{t-1}}} \left[\max_{\pi \in \Pi} Q^{\pi_{t-1}}(x, \cdot)^{\top} \pi(x) - Q^{\pi_{t-1}}(x, \cdot)^{\top} \pi'_t(x) \right] \leq \frac{2|\mathcal{A}|}{(1-\gamma)\alpha} \left(\varepsilon_0 + \frac{2}{\sqrt{N}} \right),$$

142 where α is the weak-learner multiplicative factor in [Definition 2.2](#), ε_0 is the boosting oracle error,
 143 and N is the number of weak learners used in the Internal Boost.

Algorithm 1 RL Boosting ([Brukhim et al., 2022](#))

- 1: **Input:** reward function r , policy class Π , number of iterations T , initial state distribution μ , and
 P, N, M parameters for Internal Boost.
 - 2: Initialize a policy $\pi_0 \in \Pi$ arbitrarily.
 - 3: **for** $t = 1$ **to** T **do**
 - 4: Run Internal Boost ([Algorithm 2](#) of [Brukhim et al. \(2022\)](#)) with distribution μ and policy
 π_{t-1} to obtain π'_t .
 - 5: Update $\pi_t = (1 - \eta_{1,t})\pi_{t-1} + \eta_{1,t}\pi'_t$.
 - 6: **end for**
 - 7: Run each policy π_t for P rollouts to compute an empirical estimate \widehat{V}^{π_t} of the expected return.
 - 8: **return** $\bar{\pi} := \pi_{t'}$ where $t' = \operatorname{argmax}_t \widehat{V}^{\pi_t}$.
-

144 **3 Main results**

145 In this section, we introduce our main algorithms ([Algorithm 2](#), [Algorithm 3](#)). The convergence
 146 of both algorithms depends on coverability parameters (defined in [Section 3.1](#) and [Section 3.2](#)
 147 respectively) that are known to be bounded for many interesting classes of structured MDPs. Notably,
 148 neither algorithm assumes bounded distribution mismatch C_{∞} for the underlying MDP and policy
 149 class, and so we give the first boosting algorithms with provable convergence guarantees without
 150 requiring bounded C_{∞} .

151 The assumption of bounded C_{∞} is highly restrictive. For C_{∞} to be bounded, every policy $\pi \in \Pi$
 152 must have an occupancy d^{π} that covers the optimal policy's visitation. If even a single policy in
 153 the class fails this, C_{∞} becomes unbounded. This effectively eliminates the need for exploration, a
 154 requirement that rarely holds in practice. We provide the following example showing that bounded
 155 C_{∞} is indeed a strong assumption even under a simple RL setting ([Figure 1](#)).

156 **Motivating example.** Consider the tabular MDP, where we only have one single starting state, i.e.,
 157 $d_0(x_1) = 1, d_0(x_2) = 0, d_0(x_3) = 0$. The states x_2 and x_3 are self-loop, and the agent receives zero
 158 reward when it is in either of these states. The action space consists of n actions, and transitions
 159 are deterministic. Let $\Pi = \{\pi_1, \dots, \pi_n\}$ denote the set of deterministic policies, where each policy
 160 selects a single action. The optimal policy is the deterministic policy π_n that selects a_n , which
 161 transitions to the left and yields a reward of 1. Each of the remaining deterministic policies selects
 162 one of the other actions and yields a reward of $1 - \varepsilon$. Let $\text{CH}(\Pi)$ denote the convex hull of Π , and let
 163 $\pi \in \text{CH}(\Pi)$ be any boosted (stochastic) policy. It is immediate that every policy in $\text{CH}(\Pi)$ achieves
 164 near-optimal performance: $V^* - V^\pi \leq \varepsilon, \forall \pi \in \text{CH}(\Pi)$. However, there exists a policy π in $\text{CH}(\Pi)$
 165 that assigns zero probability to the optimal action. For such a policy, the induced state distribution
 166 satisfies $d^\pi(x_2) = 0$, while $d^{\pi^*}(x_2) > 0$. As a result, $C_\infty := \sup_{\pi \in \text{CH}(\Pi)} \left\| \frac{d^{\pi^*}}{d^\pi} \right\|_\infty \rightarrow \infty$ due to
 167 the nature definition of C_∞ , a supremum over all policy in $\text{CH}(\Pi)$.

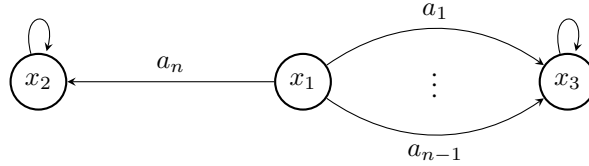


Figure 1: An MDP with 3 states and n actions, play optimal action a_n will get reward of 1, the other $n - 1$ actions will get reward of $1 - \varepsilon$, and receive no reward when in the self-loop state x_2, x_3 .

168 Our solution to eliminate the dependency on the distribution mismatch coefficient C_∞ is to incor-
 169 porate an explicit exploration strategy into the boosting framework of [Brukhim et al. \(2022\)](#) (see
 170 Appendix L.1 for a comprehensive description of their boosting algorithm, which we use as a subrou-
 171 tine). Specifically, we augment the true reward with an exploration bonus constructed via a policy
 172 cover. We leverage recent developments in L_∞ policy coverage and pushforward coverability within
 173 online reinforcement learning, which allows the algorithm to adapt automatically to the inherent
 174 complexity of the underlying MDP ([Amortila et al., 2024](#)). While L_∞ coverability generally requires
 175 stronger prior knowledge of the MDP to construct the bonus, the pushforward coverability framework
 176 relaxes this requirement by enabling the estimation of an augmented bonus.

177 Our algorithmic approach follows similarly to [Agarwal et al. \(2020\)](#), which constructs a policy
 178 cover in an iterative manner, though their results are restricted to *linear MDPs*. Using this policy
 179 cover, they design an exploration bonus that is added to the reward to encourage exploration, and
 180 optimize policies by solving $\max_{\pi} \mathbb{E}_{s_0, a_0 \sim \rho_{\text{cov}}^n} [Q^\pi(s_0, a_0; r + b)]$, where b is the bonus function
 181 and ρ_{cov} is a distribution constructed from the policy cover. In our framework, we use boosting to
 182 solve this policy optimization problem. Our theoretical guarantees are established specifically for
 183 the boosting-based solution, whereas in practice, the optimization step can be instantiated with any
 184 suitable policy optimization algorithm.

185 Next, we present our first main result on L_∞ coverability.

186 3.1 L_∞ -coverability

187 In this section, we analyze the performance of our boosting algorithm under the assumption of access
 188 to a reference distribution μ that provides global coverage of the state-action space. This ensures that
 189 any admissible occupancy distribution induced by a policy is not “too far” from μ in terms of their
 190 density ratio ([Amortila et al., 2024](#); [Chen and Jiang, 2019](#)). This requirement is formalized by the
 191 following assumption:

192 **Assumption 3.1** (L_∞ -concentrability coefficient ([Chen and Jiang, 2019](#))). We assume ac-
 193 cess to a distribution $\mu \in \Delta(\mathcal{X} \times \mathcal{A})$ for which the L_∞ -concentrability coefficient
 194 $\tilde{C}_\infty := \sup_{\pi} \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \left\{ \frac{d^\pi(x,a)}{\mu(x,a)} \right\}$ is bounded.

195 Next, we provide the pseudo-code of our first algorithm in Algorithm 2. The algorithm runs for K
 196 iterations to generate policy cover. To ensure sufficient exploration under the L_∞ -concentrability
 197 assumption, we define a weight function $w^k(x, a)$ that decays as state-action pairs are visited by

Algorithm 2 Boosting with L_∞ coverability

- 1: **input:** number of epoch K , number of iterations T , number of weak learners N , number of sample episodes M , policy class Π .
 - 2: Initialize policy π^1 arbitrarily, $p^1 = \{\pi^1\}$.
 - 3: **for** $k = 1, \dots, K$ **do**
 - 4: Construct a new bonus $b^k(x, a)$ with policy cover d^{p^k} (see Eq. (2)).
 - 5: Update policy class $\widehat{\Pi} \leftarrow \mathcal{T}(\Pi, b^k)$ (see Eq. (3)).
 - 6: Run RL boosting (Algorithm 1) with reward function $r + b^k$, policy class $\widehat{\Pi}$, initial distribution d^{p^k} to obtain π^k .
 - 7: Update the policy cover $p^{k+1} = \text{Unif}(\pi_1, \dots, \pi_k)$.
 - 8: **end for**
 - 9: **return** $\pi := \operatorname{argmax}_{\pi \in \{\pi^1, \dots, \pi^K\}} V^\pi$.
-

198 previous policies. For each iteration k , the weight function and the corresponding exploration bonus
 199 $b^k(x, a)$ are defined as follows:

$$w^k(x, a) = \frac{\widetilde{C}_\infty \mu(x, a)}{\sum_{i < k} d^{\pi_i}(x, a) + \widetilde{C}_\infty \mu(x, a)}, \quad b^k(x, a) = \frac{2}{1 - \gamma} w^k(x, a). \quad (2)$$

200 The weight $w^k(x, a)$ acts as an inverse frequency measure. When the total density of previous policies
 201 $\sum d^{\pi_i}$ at a certain (x, a) is low relative to the reference distribution μ , the weight w^k remains high,
 202 providing a larger bonus to encourage exploration. As the algorithm collects more samples and the
 203 cumulative density grows, the bonus naturally vanishes. The factor $\frac{1}{1 - \gamma}$ in the bonus is to cap the
 204 maximum value of any value function $V^\pi(x)$.

205 To facilitate exploration, the algorithm adopts a uniform strategy over actions that remain insufficiently
 206 explored. This mechanism is formalized through a policy class transformation, $\mathcal{T}(\Pi, b^k)$, which
 207 maps each original policy $\pi \in \Pi$ to a corresponding transformed policy $\widehat{\pi}$, which we formalize as
 208 follows.

209 **Policy transformation** $\mathcal{T}(\Pi, b^k)$. Given the well-explored set $\mathcal{B}^k := \{x : w^k(x, a) \leq \frac{1}{2}, \forall a \in \mathcal{A}\}$,
 210 the transformed policy $\widehat{\pi}$ is defined as follows:

$$\widehat{\pi}(x) = \begin{cases} \pi(x), & \text{if } x \in \mathcal{B}^k \\ \text{Unif}\left(\{a \in \mathcal{A} : w^k(x, a) > \frac{1}{2}\}\right), & \text{otherwise} \end{cases} \quad (3)$$

211 Intuitively, the policy simply takes a random action on the state where the bonus is already large. We
 212 note that the policy transformation is introduced purely for abstraction; in practice, we do not need to
 213 iterate over the entire policy class Π . Its implementation is straightforward by checking the condition
 214 for all the actions (e.g., Algorithm 3 in Agarwal et al. (2020)).

215 With the policy transformation in place, we present the main result for the RL boosting under L_∞
 216 coverability as follows:

Theorem 3.2. For any $\varepsilon, \delta \in (0, 1)$, suppose Assumption 3.1 is satisfied, Algorithm 2 is instantiated with the following parameters: $K = O\left(\frac{\widetilde{C}_\infty \iota}{(1 - \gamma)\varepsilon}\right)$, $T = O\left(\frac{\widetilde{C}_\infty^2 \iota^2}{(1 - \gamma)^8 \varepsilon^3}\right)$, $N = O\left(\frac{\widetilde{C}_\infty^2 |\mathcal{A}|^2 \iota^2}{(1 - \gamma)^{10} \alpha^2 \varepsilon^4}\right)$, $M = m\left(\frac{(1 - \gamma)^5 \alpha \varepsilon^2}{\widetilde{C}_\infty |\mathcal{A}|}, \frac{\delta}{KTN}\right)$, where $\iota = \log\left(\frac{\widetilde{C}_\infty}{(1 - \gamma)\varepsilon}\right)$. Then Algorithm 2 produces policy π , such that with probability at least $1 - \delta$,

$$V^* - V^\pi \leq \frac{\widetilde{C}_\infty \iota}{(1 - \gamma)^5 \varepsilon} \mathcal{E}(\Pi, \Pi) + \varepsilon.$$

217 **Sample Complexity:** If $m(\varepsilon, \delta) = \frac{\log |\mathcal{W}|}{\varepsilon^2} \log \frac{1}{\delta}$ for some measure of weak learning complexity $|\mathcal{W}|$,
 218 then the algorithm uses at most $KTNM \leq \widetilde{O}\left(\frac{\widetilde{C}_\infty^7 |\mathcal{A}|^4 \log |\mathcal{W}|}{(1 - \gamma)^{29} \alpha^4 \varepsilon^{12}}\right)$ episodes.

219 In tabular MDP, we can choose the distribution $\mu(x, a) = \frac{1}{|\mathcal{X}||\mathcal{A}|}$ which admits a closed form and
 220 $\tilde{C}_\infty(\mu) \leq |\mathcal{X}||\mathcal{A}|$. For a Block MDP with latent state space \mathcal{S} , emission distribution $q : \mathcal{S} \rightarrow \Delta(\mathcal{X})$,
 221 and decoder $\phi : \mathcal{X} \rightarrow \mathcal{S}$, the distribution $\mu(x, a) := q(x | \phi(x)) \cdot \frac{1}{|\mathcal{S}||\mathcal{A}|}$ achieves $\tilde{C}_\infty(\mu) \leq |\mathcal{S}||\mathcal{A}|$.
 222 However, the choice of distribution μ relies on prior knowledge of the underlying MDP. Pushforward
 223 coverability relaxes this requirement by constructing the exploration objective directly from the
 224 transition dynamics, which we present the corresponding result in the next section.

225 3.2 Pushforward coverability

226 The pushforward coverability is a more general notion that can be applied to various structured MDPs,
 227 and does not require the reference distribution like L_∞ coverability. The pseudo-code of the boosting
 228 with policy cover algorithm is provided in Algorithm 3. At each episode k , we use the following
 229 weight function and bonus as follows:

$$w^k(x | x', a') := \frac{P(x | x', a')}{\sum_{i < k} d^{\pi^i}(x) + P(x | x', a')}, \quad b^k(x | x', a') := \frac{2}{1 - \gamma} w^k(x | x', a'), \quad (4)$$

230 Intuitively, if a state x is poorly covered by previous policies, the sum $\sum_{i < k} d^{\pi^i}(x)$ will be small,
 231 and the resulting weight will be higher, thereby encouraging exploration of that state.

232 By using the exploration objective in Eq. (4), the sample complexity of the algorithm adapts to the
 233 intrinsic complexity of the MDP, which is captured by the pushforward coverability C_{push} , which we
 234 recall below.

Definition 3.3 (Pushforward coverability (Xie and Jiang, 2021)).

$$C_{\text{push}} = \inf_{\mu \in \Delta(\mathcal{X})} \sup_{(x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}} \left\{ \frac{P(x' | x, a)}{\mu(x')} \right\},$$

235 The pushforward coverability C_{push} is bounded by the intrinsic complexity of the MDP. This is a
 236 general notion that can be applied to various structured MDPs. For example, tabular MDPs with state
 237 space \mathcal{X} have $C_{\text{push}} \leq |\mathcal{X}|$, block MDPs with latent space \mathcal{S} have $C_{\text{push}} \leq |\mathcal{S}|$, and low-rank MDPs
 238 of dimension d have $C_{\text{push}} \leq d$.

239 Since the reward construction requires the knowledge of the transition, which is unknown in general,
 240 we can adopt the estimate the weight function via contrastive learning (Amortila et al., 2024),
 241 which we present in Algorithm 4. The weight estimation can be solved efficiently via a convex
 242 program (Nguyen et al., 2010); we refer to Katdare et al. (2023) for a practical implementation
 243 for estimating the weight. Conceptually, we estimate $\hat{w}^k(x | x', a') \approx w^k(x | x', a')$, and use
 244 $\hat{b}^k(x | x', a') := \frac{2}{1 - \gamma} \hat{w}^k(x | x', a')$ as the exploration bonus. To perform weight function estimation,
 245 we assume access to a weight function class $\overline{\mathcal{W}}$, with $\overline{\mathcal{W}} \subseteq (\mathcal{X} \times \mathcal{A} \times \mathcal{X} \rightarrow \mathbb{R}_+)$ that is capable of
 246 representing the weight function w^k in Eq. (4). We provide this assumption as follows:

Assumption 3.4 (Weight function realizability (Amortila et al., 2024)). For all $\pi \in \Pi$:

$$w(x' | x, a) := \frac{P(x' | x, a)}{d^\pi(x')} \in \overline{\mathcal{W}}.$$

247 At a high level, the estimated weight \hat{w} is closed to the true weight w under the Hellinger dis-
 248 tance: $D_{\mathbf{H}, \nu}^2(w, w') := \mathbb{E}[(\sqrt{\hat{w}} - \sqrt{w'})^2] \lesssim \frac{\|w\|_\infty \log(|\overline{\mathcal{W}}| \delta^{-1})}{n}$, where n is number of sample in the
 249 estimation procedure. Next, we present our main result for the pushforward coverability.

Theorem 3.5. For any $\varepsilon, \delta \in (0, 1)$, suppose Assumption 3.4 is satisfied. Algorithm 3 is instanti-
 ated with the following parameters: $K = O(\frac{C_{\text{push}} \iota}{(1 - \gamma)^2 \varepsilon})$, $T = O(\frac{C_{\text{push}}^2 \iota^2}{(1 - \gamma)^{18} \varepsilon^4})$, $N = O(\frac{C_{\text{push}}^2 |\mathcal{A}|^2 \iota^2}{(1 - \gamma)^{18} \alpha^2 \varepsilon^6})$,
 $M = m\left(\frac{(1 - \gamma)^9 \alpha \varepsilon^3}{C_{\text{push}} |\mathcal{A}|}, \frac{\delta}{K T N N_{\text{weight}}}\right)$, $N_{\text{weight}} = O(\frac{C_{\text{push}} |\mathcal{A}| \log(|\overline{\mathcal{W}}| \delta^{-1}) \iota}{(1 - \gamma)^{11} \varepsilon^3})$, where $\iota = \log(\frac{C_{\text{push}}}{(1 - \gamma)^2 \varepsilon})$.
 Then Algorithm 3 produces policy π , such that with probability at least $1 - \delta$,

$$V^* - V^\pi \leq \frac{C_{\text{push}} \iota}{(1 - \gamma)^6 \varepsilon} \mathcal{E}(\Pi, \Pi) + \varepsilon.$$

250 **Sample Complexity:** If $m(\varepsilon, \delta) = \frac{\log |\mathcal{W}|}{\varepsilon^2} \log \frac{1}{\delta}$ for some measure of weak learning complexity $|\mathcal{W}|$,
 251 then the algorithm uses at most $K(TNM + N_{\text{weight}}) \leq \tilde{O} \left(\frac{C_{\text{push}}^7 |\mathcal{A}|^4 \log(|\mathcal{W}|)}{(1-\gamma)^{53} \alpha^4 \varepsilon^{17}} + \frac{C_{\text{push}}^2 |\mathcal{A}| \log(|\overline{\mathcal{W}}|)}{(1-\gamma)^{13} \varepsilon^4} \right)$
 252 episodes.

253 The large polynomial dependence on $(1-\gamma)^{-1}$ arises from the fact that our algorithm operates within
 254 the policy learning regime; we do not impose any structural assumption on the Q or V function.
 255 Recent work on policy learning in the episodic setting also experienced similar scaling challenges,
 256 e.g., $\tilde{O} \left(\frac{C_{\text{push}}^4 |S|^{24} |\mathcal{A}|^{30} H^{39}}{\varepsilon^{18}} \right)$ (Krishnamurthy et al., 2025, Theorem 4), where H is the horizon.

257 Our results offer an advantage over existing boosting framework. Specifically, Brukhim et al. (2022)
 258 established an optimality gap of $V^* - V^\pi \leq \frac{C_\infty \mathcal{E}(\Pi, \Pi)}{1-\gamma} + \varepsilon$, this bound can become vacuous in the
 259 regime of unbounded C_∞ , while our bound still gives a meaningful result. In addition, the number
 260 of iterations required for their algorithm to converge depends on C_∞ , as does the number of weak
 261 learners aggregated in the final policy.

262 Algorithm 3 is oracle efficient: assuming access to an efficient weak learning oracle, and an effi-
 263 cient optimizer for the weight function, our result implies that both computational and statistical
 264 complexities scale as $\text{poly}(C_{\text{push}}, (1-\gamma)^{-1}, |\mathcal{A}|, \alpha^{-1}, \varepsilon^{-1}, \log |\mathcal{W}|, \log(|\overline{\mathcal{W}}|), \log \delta^{-1})$:

265 • *Statistical efficiency.* Theorem 3.5 shows that the sample complexity scales polynomially
 266 with the intrinsic complexity of the MDP, as captured by the pushforward coverability C_{push} ,
 267 where it is bounded in many structured MDPs.

268 • *Computational efficiency.* Algorithm 3 requires a total of $\tilde{O} \left(\frac{C_{\text{push}}^5 |\mathcal{A}|^2 \log(|\mathcal{W}|)}{(1-\gamma)^{30} \alpha^2 \varepsilon^{10}} \right)$ calls to
 269 the weak learning oracle, and K calls to an optimization oracle for the weight function (in
 270 Algorithm 4).

271 **Proof sketch.** The core of our analysis involves a novel construction of an extended MDP and a
 272 truncated policy class, building upon the framework of Amortila et al. (2024). Prior work utilizes
 273 this construction for the reward-free setting, optimizing only for the exploration bonuses. In contrast,
 274 our algorithm optimizes both the true reward and the exploration bonuses. This creates a non-trivial
 275 coupling, and poses a significant analytical challenge, particularly in the discounted MDP setting
 276 where the episodic construction of Amortila et al. (2024) does not hold. To resolve this, we propose a
 277 redesigned construction of the extended MDP and a new termination condition of the truncated policy
 278 class. This construction allows us to partition the state space into well-covered and poorly covered
 279 regions, effectively managing the latter via our exploration bonus and the Performance Difference
 280 Lemma. Consequently, we can bound the distribution mismatch in the well-covered regions by
 281 $\left\| \frac{d^{\pi^*}}{(1-\gamma)d^{p^k}} \right\|_\infty$ with d^{p^k} as our exploration distribution induced by the policy cover. Crucially, this
 282 yields a gradient dominance property that is independent of the global distribution mismatch, enabling
 283 the direct application of non-convex Frank-Wolfe convergence guarantees of Brukhim et al. (2022).
 284 Furthermore, the pushforward analysis requires exploration bonuses to reach a $\frac{1}{1-\gamma}$ threshold on
 285 poorly explored states: a condition that typically relies on a known lower bound on MDP transition
 286 dynamics. We relax this requirement through an adaptive tuning technique: by treating the transition
 287 threshold as a variable within the truncated policy condition and optimizing it in the final bound, we
 288 eliminate the need for prior knowledge of the lower bound on the transition dynamics. Importantly,
 289 these constructions — the extended MDP, truncated policies and adaptive threshold — serve purely as
 290 analytical tools; the underlying algorithms remain agnostic to them, preserving its practical simplicity.

Algorithm 3 Boosting with push forward coverability

- 1: **input:** total iterations T , number of weak learners N , total sample episodes M , policy class Π .
 - 2: Initialize policy π^1 arbitrarily, $p^1 = \{\pi^1\}$.
 - 3: **for** $k = 1, \dots, K$ **do**
 - 4: Run Algorithm 4 to estimate $\widehat{w}^k(x | x', a')$ with cover p^k and define bonus $\widehat{b}^k(x | x', a')$.
 - 5: Run RL boosting (Algorithm 1) with reward $r + \widehat{b}^k$, initial distribution d^{p^k} to obtain π^k .
 - 6: Update the policy cover $p^{k+1} = \text{Unif}(\pi_1, \dots, \pi_k)$.
 - 7: **end for**
 - 8: **return** $\pi := \operatorname{argmax}_{\pi \in \{\pi^1, \dots, \pi^K\}} V^\pi$.
-

Algorithm 4 Weight estimation

- 1: **function** SAMPLECOVER
 - 2: Sample $x_0 \sim \mu$, continue to execute π , at any step h with (x_{h-1}, a_{h-1}, x_h) terminate with probability $1 - \gamma$.
 - 3: Play action $a \sim \text{Unif}(\mathcal{A})$, and observe (x_h, a, x_{h+1}) .
 - 4: **return** (x_h, a, x_{h+1}) .
 - 5: **end function**
 - 6: **function** SAMPLESTATEACTION
 - 7: Sample $x_0 \sim \mu$, continue to execute π , at any step h with (x_{h-1}, a_{h-1}, x_h) terminate with probability $1 - \gamma$.
 - 8: **return** (x_{h-1}, a_{h-1}, x_h) .
 - 9: **end function**
 - 10: **input:** policy cover p .
 - 11: Initialize $\mathcal{D}_1 = \mathcal{D}_2 = \emptyset$.
 - 12: For each $j \in [n]$, draw $\pi \sim p$ and sample $(x'_j, a'_j, x_j) \sim \pi$ with SampleCover. Add $(x'_j, a'_j, x_j) \sim \pi$ to both \mathcal{D}_1 and \mathcal{D}_2 .
 - 13: **for** $i < k$ **do**
 - 14: Draw $\{(x'_j, a'_j, x_j)\}_{j \in [n]}$ by drawing $\pi \sim p$ and $(x'_j, a'_j, x_j) \sim \pi$ with SampleCover.
 - 15: Draw $\{(\tilde{x}'_j, \tilde{a}'_j, \tilde{x}_j)\}_{j \in [n]}$ by sampling $(\tilde{x}'^j, \tilde{a}'^j, \tilde{x}^j) \sim \pi^i$ with SampleStateAction.
 - 16: Add $\{(x'_j, a'_j, x_j)\}_{j \in [n]}$ to \mathcal{D}_1 and add $\{(x'_j, a'_j, \tilde{x}_j)\}_{j \in [n]}$ to \mathcal{D}_2 .
 - 17: **end for**
 - 18: Set $\widehat{w} := \operatorname{argmax}_{w \in \overline{\mathcal{W}}} \widehat{\mathbb{E}}_{\mathcal{D}_1} \left[\log \left(w(x_h | x_{h-1}, a_{h-1}) \right) \right] - k \cdot \widehat{\mathbb{E}}_{\mathcal{D}_2} \left[w(x_h | x_{h-1}, a_{h-1}) \right]$.
-

291 **4 Conclusions**

292 In this work, we propose a boosting-based approach for reinforcement learning in structured MDPs.
293 Rather than relying on distribution mismatch assumptions commonly used in prior work, we use
294 an intrinsic notion of coverability, which is bounded in a wide class of structured MDPs. Under
295 this condition, we show that the proposed boosting method is both statistically efficient and oracle
296 efficient, with complexity scaling polynomial in the relevant problem parameters.

297 Several promising directions remain for future research. Our approach leverages policy covers to guide
298 exploration in a principled manner, which plays a crucial role in relaxing dependence on distribution
299 mismatch. It would be interesting to investigate alternative mechanisms for strategic exploration
300 that could further broaden the applicability of the framework while maintaining computational and
301 statistical efficiency. The guarantee provided in this paper assumes access to a *multiplicative* weak-
302 learner, further development for other types of weak learners is another natural direction. While
303 our algorithms have complexities that depend only polynomially on all relevant parameters, the
304 large polynomial dependence on $1 - \gamma$ is undesirable, and we leave open whether this and other
305 dependencies can be improved with alternative approaches. Furthermore, our results show that,
306 given access to a weak learner, efficient boosting is achievable under the notion of pushforward
307 coverability. Whether such a coverability condition is also necessary for efficient boosting remains
308 an open question for future work.

309 **References**

- 310 Alekh Agarwal, Mikael Henaff, Sham Kakade, and Wen Sun. Pc-pg: Policy cover directed exploration
311 for provable policy gradient learning. *Advances in neural information processing systems*, 33:
312 13399–13412, 2020.
- 313 Naman Agarwal, Brian Bullins, and Karan Singh. Variance-reduced conservative policy iteration. In
314 Shipra Agrawal and Francesco Orabona, editors, *Proceedings of The 34th International Conference*
315 *on Algorithmic Learning Theory*, volume 201 of *Proceedings of Machine Learning Research*,
316 pages 3–33. PMLR, 20 Feb–23 Feb 2023. URL [https://proceedings.mlr.press/v201/
317 agarwal23a.html](https://proceedings.mlr.press/v201/agarwal23a.html).
- 318 Philip Amortila, Dylan J Foster, and Akshay Krishnamurthy. Scalable online exploration via
319 coverability. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver,
320 Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference*
321 *on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 1392–
322 1455. PMLR, 21–27 Jul 2024. URL [https://proceedings.mlr.press/v235/amortila24a.
323 html](https://proceedings.mlr.press/v235/amortila24a.html).
- 324 Ronen I Brafman and Moshe Tennenholtz. R-max-a general polynomial time algorithm for near-
325 optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231, 2002.
- 326 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and
327 Wojciech Zaremba. Openai gym, 2016.
- 328 Nataly Brukhim, Elad Hazan, and Karan Singh. A boosting approach to reinforcement learning.
329 *Advances in Neural Information Processing Systems*, 35:33806–33817, 2022.
- 330 Kai-Wei Chang, Akshay Krishnamurthy, Alekh Agarwal, Hal Daumé III, and John Langford. Learning
331 to search better than your teacher. In *International Conference on Machine Learning*, pages 2058–
332 2066. PMLR, 2015.
- 333 Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In
334 *International conference on machine learning*, pages 1042–1051. PMLR, 2019.
- 335 Ching-An Cheng, Andrey Kolobov, and Alekh Agarwal. Policy improvement via imitation of multiple
336 oracles. *Advances in Neural Information Processing Systems*, 33:5587–5598, 2020.
- 337 Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford.
338 Provably efficient rl with rich observations via latent state decoding. In *International Conference*
339 *on Machine Learning*, pages 1665–1674. PMLR, 2019.
- 340 John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra. Efficient projections onto the
341 l_1 -ball for learning in high dimensions. In *Proceedings of the 25th international conference on*
342 *Machine learning*, pages 272–279, 2008.
- 343 Y. Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121
344 (2):256–285, 1995. ISSN 0890-5401. doi: <https://doi.org/10.1006/inco.1995.1136>. URL [https:
345 //www.sciencedirect.com/science/article/pii/S0890540185711364](https://www.sciencedirect.com/science/article/pii/S0890540185711364).
- 346 Udaya Ghai and Karan Singh. Sample-optimal agnostic boosting with unlabeled data. In *Forty-second*
347 *International Conference on Machine Learning*, 2025. URL [https://openreview.net/forum?
348 id=hcLeFe7idT](https://openreview.net/forum?id=hcLeFe7idT).
- 349 Elad Hazan and Karan Singh. Boosting for online convex optimization. In *International Conference*
350 *on Machine Learning*, pages 4140–4149. PMLR, 2021.
- 351 Marcel Hussing, Michael Kearns, Aaron Roth, Sikata B Sengupta, and Jessica Sorrell. Oracle-efficient
352 reinforcement learning for max value ensembles. *Advances in Neural Information Processing*
353 *Systems*, 37:117657–117681, 2024.
- 354 Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Context-
355 tual decision processes with low bellman rank are pac-learnable. In *International Conference on*
356 *Machine Learning*, pages 1704–1713. PMLR, 2017.

- 357 Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient?
358 *Advances in neural information processing systems*, 31, 2018.
- 359 Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement
360 learning with linear function approximation. In *Conference on learning theory*, pages 2137–2143.
361 PMLR, 2020.
- 362 Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In
363 *Proceedings of the nineteenth international conference on machine learning*, pages 267–274, 2002.
- 364 Daniel Kane, Sihan Liu, Shachar Lovett, and Gaurav Mahajan. Computational-statistical gap in
365 reinforcement learning. In *Conference on Learning Theory*, pages 1282–1302. PMLR, 2022.
- 366 Pulkit Katdare, Nan Jiang, and Katherine Rose Driggs-Campbell. Marginalized importance sampling
367 for off-environment policy evaluation. In *Conference on Robot Learning*, pages 3778–3788. PMLR,
368 2023.
- 369 Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time.
370 *Machine learning*, 49:209–232, 2002.
- 371 Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Pac reinforcement learning with rich
372 observations. *Advances in Neural Information Processing Systems*, 29, 2016.
- 373 Akshay Krishnamurthy, Gene Li, and Ayush Sekhari. The role of environment access in agnostic
374 reinforcement learning. *arXiv preprint arXiv:2504.05405*, 2025.
- 375 Andrey Kurenkov, Ajay Mandlekar, Roberto Martin-Martin, Silvio Savarese, and Animesh Garg.
376 Ac-teach: A bayesian actor-critic method for policy learning with an ensemble of suboptimal
377 teachers. *arXiv preprint arXiv:1909.04121*, 2019.
- 378 Guohao Li, Matthias Mueller, Vincent Casser, Neil Smith, Dominik L Michels, and Bernard Ghanem.
379 Oil: Observational imitation learning. *arXiv preprint arXiv:1803.01129*, 2018.
- 380 Xuefeng Liu, Takuma Yoneda, Chaoqi Wang, Matthew Walter, and Yuxin Chen. Active policy
381 improvement from multiple black-box oracles. In *International Conference on Machine Learning*,
382 pages 22320–22337. PMLR, 2023.
- 383 Xuefeng Liu, Takuma Yoneda, Rick Stevens, Matthew Walter, and Yuxin Chen. Blending imitation
384 and reinforcement learning for robust policy improvement. In *The Twelfth International Conference*
385 *on Learning Representations*, 2024. URL <https://openreview.net/forum?id=eJ0dzPJq1F>.
- 386 Zak Mhammedi, Adam Block, Dylan J Foster, and Alexander Rakhlin. Efficient model-free explo-
387 ration in low-rank mdps. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine,
388 editors, *Advances in Neural Information Processing Systems*, volume 36, pages 66782–66817.
389 Curran Associates, Inc., 2023a. URL [https://proceedings.neurips.cc/paper_files/
390 paper/2023/file/d2dc4d6c7b102d05f111c02a32e7c6bc-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/d2dc4d6c7b102d05f111c02a32e7c6bc-Paper-Conference.pdf).
- 391 Zakaria Mhammedi, Dylan J Foster, and Alexander Rakhlin. Representation learning with multi-step
392 inverse kinematics: An efficient and optimal approach to rich-observation rl. In *International*
393 *Conference on Machine Learning*, pages 24659–24700. PMLR, 2023b.
- 394 Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstrac-
395 tion and provably efficient rich-observation reinforcement learning. In *International conference on*
396 *machine learning*, pages 6961–6971. PMLR, 2020.
- 397 XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals
398 and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*,
399 56(11):5847–5861, 2010.
- 400 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
401 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
402 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
403 12:2825–2830, 2011.

- 404 Robert E Schapire. The strength of weak learnability. *Machine learning*, 5(2):197–227, 1990.
- 405 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional
406 continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*,
407 2015.
- 408 Yanjie Song, Ponnuthurai Nagarathnam Suganthan, Witold Pedrycz, Junwei Ou, Yongming He,
409 Yingwu Chen, and Yutong Wu. Ensemble reinforcement learning: A survey. *Applied Soft
410 Computing*, page 110975, 2023.
- 411 Wen Sun, Arun Venkatraman, Geoffrey J Gordon, Byron Boots, and J Andrew Bagnell. Deeply
412 aggravated: Differentiable imitation learning for sequential prediction. In *International conference
413 on machine learning*, pages 3309–3318. PMLR, 2017.
- 414 Tengyang Xie and Nan Jiang. Batch value-function approximation with only realizability. In Marina
415 Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine
416 Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 11404–11413. PMLR,
417 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/xie21d.html>.
- 418 Tengyang Xie, Dylan J Foster, Yu Bai, Nan Jiang, and Sham M Kakade. The role of coverage in
419 online reinforcement learning. In *ICLR*, 2023.

420 **A Notation: List of Symbols**

421 **Markov Decision Process**

\mathcal{X}	Observation space
\mathcal{S}	Latent state space
\mathcal{A}	Action space
$\Delta_{\mathcal{A}}$	Probability simplex over actions
$Q^{\pi}(s, a)$	Q function
$V^{\pi}(s)$	Value function
422 $d^{\pi}(x)$	state distribution
γ	Discount factor
$\mathcal{E}_{\mu}(\mathbb{P}, \Pi)$	Policy completeness
C_{∞}	Distribution mismatch
\tilde{C}_{∞}	L_{∞} coefficient
C_{push}	pushforward coverability

423 **Weak Learning and Boosting**

α	Weak learning parameter
424 $\Gamma[\cdot]$	Policy projection
Π	Policy class
\mathbb{P}	Policy-Tree class (w.r.t Π)

425 **B Related work**

426 The boosting approach is a prominent strategy for provably reducing the task of strong learning
427 to that of weak learning of classifiers that are simpler to compute but only slightly better than
428 random guessing. While boosting is a well-established technique in supervised learning with strong
429 theoretical guarantees (Schapire, 1990; Freund, 1995), its adaptation to reinforcement learning (RL)
430 has emerged only recently (Brukhim et al., 2022; Agarwal et al., 2023; Ghai and Singh, 2025). These
431 boosting frameworks can achieve global convergence guarantees under the assumption of access to a
432 multiplicative weak learner. However, their analysis crucially depends on occupancy mismatch terms
433 relative to the optimal policy. These terms can grow unbounded unless the policy class provides
434 sufficient state-space coverage, highlighting a key limitation of their approach.

435 The issue of occupancy mismatch suggests that strategic exploration is necessary to ensure that the
436 policy class adequately covers the state space. A large body of work in value-based RL has studied
437 the optimism principle, which introduces an exploration bonus to encourage visiting underexplored
438 states. Such bonuses promote exploration by increasing the estimated value of rarely visited states, as
439 in Upper Confidence Bound (UCB) methods (Jin et al., 2018, 2020). An alternative line of research
440 focuses on reward-free exploration, where the goal is to construct a policy cover that ensures sufficient
441 reachability across states, independent of specific reward functions. Recent works have explored
442 various notions of coverage and their implications for efficient exploration in RL (Misra et al.,
443 2020; Du et al., 2019; Mhammedi et al., 2023b,a). More recently, Amortila et al. (2024) proposed
444 a generalized exploration objective that unifies prior schemes while enabling intrinsic complexity
445 control, efficient planning, and scalable exploration, and we rely heavily on this objective in our work.

446 Beyond boosting, more general ensembling methods are an active area of research in RL. Given a pre-
447 trained collection of constituent policies and a squared error regression oracle for the relevant value
448 function class, Hussing et al. (2024) shows how to ensemble the constituent policies to be competitive
449 with a max-following policy (the policy which, in state s , acts consistently with whichever constituent

450 policy has the highest value in state s). Cheng et al. (2020); Liu et al. (2023, 2024) use policy gradient
 451 methods to compete with a similar, but incomparable benchmark class of max-following policies with
 452 one-step look-ahead. Notably, these works all consider the question of ensembling non-adaptively;
 453 all constituent policies are trained up front, in comparison to boosting, in which the performance of
 454 the current policy affects the selection of the next policy to be included in the ensemble. We refer the
 455 reader to Song et al. (2023) for a survey of ensembling techniques.

456 On the more empirical side, (Li et al., 2018; Kurenkov et al., 2019) study the problem of imitation
 457 learning from multiple policies and Sun et al. (2017); Schulman et al. (2015); Chang et al. (2015)
 458 give algorithms for iterative policy improvement, but these methods are heuristic and lack strong
 459 provable guarantees of performance.

460 C Extended MDP and truncated policy class setup

461 **Extended MDP $\bar{\mathcal{M}}$.** We define an extended MDP $\bar{\mathcal{M}}$ by extending the state and action spaces such
 462 that $\bar{\mathcal{A}} = \mathcal{A} \cup \{\mathfrak{t}\}$ and $\bar{\mathcal{X}} = \mathcal{X} \cup \{\mathfrak{t}\}$. The extended MDP $\bar{\mathcal{M}}$ matches the original MDP \mathcal{M} , with the
 463 follow additions:

- 464 • **Termination:** Taking the terminal action \mathfrak{t} from any state $x \in \mathcal{X}$ results in a deterministic
 465 transition to \mathfrak{t} , and \mathfrak{t} is a self-looping terminal state.
- 466 • **Absorbing State:** The terminal state \mathfrak{t} is a self-looping absorbing state, where $P(\mathfrak{t} | \mathfrak{t}, a) = 1$
 467 for all $a \in \bar{\mathcal{A}}$.
- 468 • **Reward Structure:** Entering the terminal state via action \mathfrak{t} yields an immediate reward
 469 $r(x, \mathfrak{t}) = \frac{1}{1-\gamma}$. Once in the terminal state, all subsequent rewards are zero.

470 **Truncated policy class $\bar{\Pi}$.** Let $\mathcal{B} \subseteq \mathcal{X}$ denote a “well-explored” set (to be formally defined based
 471 on specific coverability notions). For every policy π in the original class Π , we define a corresponding
 472 truncated policy $\bar{\pi} \in \bar{\Pi}$ as:

$$\bar{\pi}(x) = \begin{cases} \pi(x), & x \in \mathcal{B} \\ \mathfrak{t}, & \text{otherwise,} \end{cases}$$

473 We note that the extended MDP and truncated policy class are defined only for the purpose of analysis.
 474 The algorithm itself does not require knowledge of these extensions, nor does it explicitly utilize the
 475 terminal state \mathfrak{t} or the set \mathcal{B} during execution.

476 **Notation.** We use the bar notation (e.g., \bar{f}) to denote quantities related to the extended MDP and
 477 the truncated policy class. Specifically, for each episode $k \in [K]$ let $\bar{\Pi}_k$ denote the truncated policy
 478 class constructed from Π using \mathcal{B}^k . For any $\pi \in \Pi$ we let $\bar{\pi} \in \bar{\Pi}_k$ be the correspond truncated policy.
 479 Finally, we denote the state and state action value functions within the extended MDP as $\bar{V}(x)$ and
 480 $\bar{Q}(x, a)$, respectively.

481 We provide the following claim for *extended* MDP with a bonus function.

482 **Claim 1.** Let $\bar{V}_b^{\bar{\pi}}$ denote the value function of the truncated policy $\bar{\pi}$ of π on the extended MDP with
 483 a positive bonus function b . Then, for any $x \in \mathcal{X}$, we have

$$V^\pi(x) \leq \bar{V}_b^{\bar{\pi}}(x),$$

484 *Proof.* **In-set transitions:** As long as the state x remains within \mathcal{B} , the truncated policy $\bar{\pi}$ follows
 485 the original policy π . In this regime, the reward in the extended MDP is greater than the original
 486 reward because the bonus is positive.

487 **Out-of-set transitions:** Upon encountering a state $x \notin \mathcal{B}$, the policy $\bar{\pi}$ takes the termination action \mathfrak{t} .
 488 This yields an immediate value of $\frac{1}{1-\gamma}$. Since the maximum possible value in the original MDP is
 489 bounded by $\frac{1}{1-\gamma}$, the value obtained by terminating is at least as large as the value of continuing with
 490 the original policy π .

491 Combining these cases, the value function in the extended MDP serves as an upper bound for the
 492 original value function. \square

493 **D Proof of Section 3.1**

Theorem 3.2. For any $\varepsilon, \delta \in (0, 1)$, suppose Assumption 3.1 is satisfied, Algorithm 2 is instantiated with the following parameters: $K = O(\frac{\tilde{C}_{\infty} \iota}{(1-\gamma)\varepsilon})$, $T = O(\frac{\tilde{C}_{\infty}^2 \iota^2}{(1-\gamma)^8 \varepsilon^3})$, $N = O(\frac{\tilde{C}_{\infty}^2 |\mathcal{A}|^2 \iota^2}{(1-\gamma)^{10} \alpha^2 \varepsilon^4})$, $M = m\left(\frac{(1-\gamma)^5 \alpha \varepsilon^2}{\tilde{C}_{\infty} |\mathcal{A}|}, \frac{\delta}{KTN}\right)$, where $\iota = \log(\frac{\tilde{C}_{\infty}}{(1-\gamma)\varepsilon})$. Then Algorithm 2 produces policy π , such that with probability at least $1 - \delta$,

$$V^* - V^{\pi} \leq \frac{\tilde{C}_{\infty} \iota}{(1-\gamma)^5 \varepsilon} \mathcal{E}(\Pi, \Pi) + \varepsilon.$$

494 **Sample Complexity:** If $m(\varepsilon, \delta) = \frac{\log |\mathcal{W}|}{\varepsilon^2} \log \frac{1}{\delta}$ for some measure of weak learning complexity $|\mathcal{W}|$,
 495 then the algorithm uses at most $KTNM \leq \tilde{O}\left(\frac{\tilde{C}_{\infty}^7 |\mathcal{A}|^4 \log |\mathcal{W}|}{(1-\gamma)^{29} \alpha^4 \varepsilon^{12}}\right)$ episodes.

496 *Proof of Theorem 3.2.* For any episode $k \in [K]$, we define $\mathcal{B}^k := \{x, a : \frac{\tilde{C}_{\infty} \mu(x, a)}{d^{p^k}(x, a)} \leq k - 1\}$ to
 497 be the set of a well-explored states relative to the reference distribution μ . We say $x \in \mathcal{B}^k$ if
 498 $(x, a) \in \mathcal{B}^k, \forall a \in \mathcal{A}$.

499 Let $\pi = \operatorname{argmin}_{\pi \in \{\pi^1, \dots, \pi^K\}} V^{\pi^k}$ be the output of the Algorithm 2, we have the following:

$$V_{d_0}^{\pi^*} - V_{d_0}^{\pi} \leq \frac{1}{K} \sum_{k=1}^K V_{d_0}^{\pi^*} - V_{d_0}^{\pi^k} \leq \frac{1}{K(1-\gamma)} \sum_{k=1}^K V_{d^{p^k}}^{\pi^*} - V_{d^{p^k}}^{\pi}. \quad (5)$$

500 In the last inequality, we use the following properties to switch to using d^{p^k} as our initial distribution
 501 in each round, for which we pay a factor of $\frac{1}{1-\gamma}$.

$$\begin{aligned} V_{d_0}^{\pi^*} - V_{d_0}^{\pi} &= \mathbb{E}_{x \sim d_0} \left[V^{\pi^*}(x) - V^{\pi}(x) \right] \\ &= \mathbb{E}_{x \sim d^{p^k}} \left[\frac{d_0(x)}{d^{p^k}(x)} (V^{\pi^*}(x) - V^{\pi}(x)) \right] \\ &\leq \frac{1}{1-\gamma} \left(V_{d^{p^k}}^{\pi^*} - V_{d^{p^k}}^{\pi} \right), \end{aligned}$$

502 where the last inequality follows from the fact that, $V^*(x) - V^{\pi}(x) \geq 0$ and
 503 $d^{p^k}(x) \geq (1-\gamma)d_0(x), \forall x \in \mathcal{X}$.

504 To reduce notational clutter throughout the remainder of the proof, we assume d^{p^k} is the initial
 505 distribution and omit the subscript d^{p^k} when referring to the value functions, e.g., we simply write
 506 $V_{d^{p^k}}^{\pi}$ as V^{π} .

507 We proceed to bound the RHS of Eq. (5) by decomposing the summation into the following three
 508 terms:

$$\frac{1}{K} \sum_{k=1}^K V^{\pi^*} - V^{\pi^k} = \underbrace{\frac{1}{K} \sum_{k=1}^K V^{\pi^*} - \bar{V}_{b^k}^{\pi^*}}_{I_1} + \underbrace{\frac{1}{K} \sum_{k=1}^K \bar{V}_{b^k}^{\pi^*} - \bar{V}_{b^k}^{\pi^k}}_{I_2} + \underbrace{\frac{1}{K} \sum_{k=1}^K \bar{V}_{b^k}^{\pi^k} - V^{\pi^k}}_{I_3},$$

509 where $\bar{\pi}^*$ be the truncated policy counterpart of π^* . First, we have $I_1 \leq 0$ follows by Claim 1.

510 For I_3 , we have

$$\begin{aligned}
\sum_{k=1}^K \bar{V}_{b^k}^{\pi^k} - V^{\pi^k} &= \sum_{k=1}^K V_{b^k}^{\pi^k} - V^{\pi^k} \\
&= 2 \sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\frac{\tilde{C}_\infty \mu(x, a)}{\sum_{i < k} d^{\pi^i}(x, a) + \tilde{C}_\infty \mu(x, a)} \right] \\
&= 2 \tilde{C}_\infty \sum_{x \in \mathcal{X}, a \in \mathcal{A}} \sum_{k=1}^K \mu(x, a) \frac{d^{\pi^k}(x, a)}{\sum_{i < k} d^{\pi^i}(x, a) + \tilde{C}_\infty \mu(x, a)}
\end{aligned}$$

511 Since $\sup_{\pi \in \Pi} d^\pi(x, a) \leq \tilde{C}_\infty \mu(x, a)$ for all $(x, a) \in \mathcal{X} \times \mathcal{A}$, Lemma L.3 implies that

$$\sum_{x \in \mathcal{X}, a \in \mathcal{A}} \sum_{k=1}^K \mu(x, a) \frac{d^{\pi^k}(x, a)}{\sum_{i < k} d^{\pi^i}(x, a) + \tilde{C}_\infty \mu(x, a)} \leq 2 \log(2K)$$

512 To bound I_2 , we proceed to bridge the *extended* bonus MDP to the result of the bonus MDP in order
513 to apply Lemma L.10 for non-convex Frank-Wolfe optimization in Appendix L.1, which requires the
514 guarantee on the gradient domination which we provide as follows:

515 **Lemma D.1.** For any $k \in [K]$, let $\kappa = \frac{k-1}{1-\gamma}$, $\tau = \frac{2(k-1)}{(1-\gamma)^3} \mathcal{E}(\Pi, \Pi)$. Let $\bar{\pi} \in \bar{\Pi}_k$ be a policy in the
516 truncated class, let $\hat{\Pi} := \pi \in \mathcal{T}(\Pi, b^k)$ be the transformed policy class, and $\hat{\Pi} := \mathcal{T}(\Pi, b^k)$ be the
517 policy tree constructed from $\hat{\Pi}$. Consider the value function on the extended bonus MDP under the
518 initial distribution d^{p^k} , we have $\bar{V}_{b^k}^{\bar{\pi}}$ is (κ, τ, Π, Π) -gradient dominated. That is, for any $\pi \in \Pi$,

$$\max_{\bar{\pi} \in \bar{\Pi}_k} \bar{V}_{b^k}^{\bar{\pi}} - \bar{V}_{b^k}^\pi \leq \frac{k-1}{1-\gamma} \max_{\pi' \in \Pi} (\nabla \bar{V}_{b^k}^\pi)^\top (\pi' - \pi) + \frac{2(k-1)}{(1-\gamma)^3} \mathcal{E}(\Pi, \Pi).$$

519 Since $\bar{Q}_{b^k} = Q_{b^k}, \forall x \in \Pi$ since all policy in Π do not pick t , and by Lemma L.7, the Internal Boost
520 guarantee that

$$\max_{\pi \in \Pi} \mathbb{E}_{(s, \bar{Q}) \sim \mathcal{D}_t} \left[\bar{Q}^\top \pi(s) \right] - \mathbb{E}_{(s, \bar{Q}) \sim \mathcal{D}_t} \left[\bar{Q}^\top \pi_t(s) \right] \leq \frac{2}{(1-\gamma)^2 \alpha} (\varepsilon_w + 2/\sqrt{N}) \quad (6)$$

521 By Eq. (6), Lemma D.1, and $\bar{V}_{b^k}^\pi$ is $\frac{4\gamma}{(1-\gamma)^3}$ -smoothness, we can appeal to Lemma L.10, put
522 everything together, we have

$$V_{d_0}^{\pi^*} - V_{d_0}^{\pi^{k'}} \leq \frac{2K}{(1-\gamma)^4} \mathcal{E}(\Pi, \Pi) + \frac{8K^2}{(1-\gamma)^6 T} + \frac{4\varepsilon_w K |\mathcal{A}|}{(1-\gamma)^4 \alpha} + \frac{8K |\mathcal{A}|}{(1-\gamma)^4 \alpha \sqrt{N}} + \frac{1}{K(1-\gamma)} \tilde{C}_\infty \log(2K),$$

523 which also concludes the proof.

524

□

525 D.1 Support proof of L_∞ coverability

526 **Claim 2.** For any policy $\bar{\pi} \in \bar{\Pi}_k$, and $\pi \in \mathcal{T}(\Pi, \mathcal{B}^k)$, for any $x \in \mathcal{B}^k$, we have

$$d^{\bar{\pi}}(x) \leq \beta^k d^{p^k}(x)$$

527 for any $x \notin \mathcal{B}^k$, we have

$$\bar{Q}_{b^k}^{\bar{\pi}}(x, t) \leq \bar{V}_{b^k}^\pi(x).$$

528 *Proof of Claim 2. Part 1.* For any $x \in \mathcal{B}^k$, $\bar{\pi}$ does not pick the terminal action, so we only need to
529 consider the original action set:

$$d^\pi(x) = \sum_{a \in \mathcal{A}} d^\pi(x, a) \leq \sum_{a \in \mathcal{A}} \tilde{C}_\infty \mu(x, a) \leq \sum_{a \in \mathcal{A}} \beta^k d^{p^k}(x, a) \leq \beta^k d^{p^k}(x),$$

530 The first inequality follows by the definition of \widetilde{C}_∞ , the second inequality due to $\frac{\widetilde{C}_\infty \mu(x,a)}{d^{p^k}(x,a)} \leq$
 531 $\beta^k, \forall a \in \mathcal{A}$ if $x \in \mathcal{B}^k$.

532 **Part 2.** By the definition of the truncated policy, for any $x \notin \mathcal{B}^k$, the policy deterministically selects
 533 the terminal action, i.e., $\pi(x) = \mathfrak{t}$. Choosing action \mathfrak{t} yields

$$\overline{Q}_{b^k}^\pi(x, \mathfrak{t}) = 1 + \gamma \overline{V}_{b^k}^\pi(x)$$

534 since the transition under \mathfrak{t} deterministically returns to state x . Moreover, because π selects actions
 535 uniformly from the set $\{a : (x, a) \notin \mathcal{B}^k\}$, implying $\overline{V}_{b^k}^\pi(x) \geq \frac{1}{1-\gamma}$. Therefore,

$$\overline{Q}_{b^k}^\pi(x, \mathfrak{t}) - \overline{V}_{b^k}^\pi(x) = 1 - (1-\gamma)\overline{V}_{b^k}^\pi(x) \leq 0$$

536 This completes the proof. \square

537 *Proof of Lemma D.1.* Let $\bar{\pi}$ be any truncated policy, by Claim 2, we have $\overline{Q}_{b^k}^{\bar{\pi}}(x, \mathfrak{t}) \leq \overline{V}_{b^k}^\pi(x)$ for
 538 all $x \notin \mathcal{B}^k$, apply Lemma L.1 yielding

$$\begin{aligned} \overline{V}_{b^k}^{\bar{\pi}} - \overline{V}_{b^k}^\pi &\leq \frac{1}{1-\gamma} \sum_{x,a} d^{\bar{\pi}}(x) \bar{\pi}(a|x) \left(\overline{Q}_{b^k}^{\bar{\pi}}(x, a) - \overline{V}_{b^k}^\pi(x) \right) \mathbb{1}\{x \in \mathcal{B}^k\} \\ &\leq \frac{k-1}{(1-\gamma)^2} \sum_{x,a} d^\pi(x) \bar{\pi}(a|x) \left(\overline{Q}_{b^k}^\pi(x, a) - \overline{V}_{b^k}^\pi(x) \right) \mathbb{1}\{x \in \mathcal{B}^k\} \\ &\leq \frac{k-1}{(1-\gamma)^2} \sum_x d^\pi(x) \left(\max_a \overline{Q}_{b^k}^\pi(x, a) - \overline{V}_{b^k}^\pi(x) \right) \\ &= \frac{k-1}{(1-\gamma)^2} \max_{\pi_0} \sum_{x,a} d^\pi(x) \pi_0(a|x) \left(\overline{Q}_{b^k}^\pi(x, a) - \overline{V}_{b^k}^\pi(x) \right) \\ &= \frac{k-1}{1-\gamma} \max_{\pi_0} (\nabla \overline{V}_{b^k}^\pi)^\top (\pi_0 - \pi), \end{aligned}$$

539 where the second inequality follows by $\frac{d^{\bar{\pi}}(x)}{d^\pi(x)} \leq \frac{d^{\bar{\pi}}(x)}{d^{p^k}(x)(1-\gamma)} \leq \frac{k-1}{1-\gamma}$.

540 For any $\pi \in \Pi$, we decompose

$$\max_{\pi_0} (\nabla \overline{V}_{b^k}^\pi)^\top (\pi_0 - \pi) \leq \left(\max_{\pi_0} (\nabla \overline{V}_{b^k}^\pi)^\top \pi_0 - \max_{\pi' \in \Pi} (\nabla \overline{V}_{b^k}^\pi)^\top \pi' \right) + \max_{\pi' \in \Pi} (\nabla \overline{V}_{b^k}^\pi)^\top (\pi' - \pi).$$

541 Finally,

$$\begin{aligned} \max_{\pi_0} (\nabla \overline{V}_{b^k}^\pi)^\top \pi_0 - \max_{\pi' \in \Pi} (\nabla \overline{V}_{b^k}^\pi)^\top \pi' &= \min_{\pi' \in \Pi} \frac{1}{1-\gamma} \mathbb{E}_{x \sim d^\pi} \left[\max_a \overline{Q}_{b^k}^\pi(x, a) - \overline{Q}_{b^k}^\pi(x, \cdot)^\top \pi' \right] \\ &\leq \frac{2}{(1-\gamma)^2} \mathcal{E}(\Pi, \Pi), \end{aligned}$$

542 where the additional factor $2/(1-\gamma)$ follows from the bound $r(x, a) + b^k(x, a) \leq 2/(1-\gamma)$ and
 543 the definition of $\mathcal{E}(\Pi, \Pi)$. Combining the above inequalities completes the proof. \square

544 E Proof for Section 3.2

545 For any episode $k \in [K]$, we define a well-explored set $\mathcal{B}^k := \{x : \frac{1}{d^{p^k}(x)} \leq \frac{k-1}{\zeta}\}$, where ζ is a
 546 parameter will be tuned later in the subsequent analysis.

Theorem 3.5. *For any $\varepsilon, \delta \in (0, 1)$, suppose Assumption 3.4 is satisfied. Algorithm 3 is instantiated with the following parameters: $K = O(\frac{C_{\text{push}} \iota}{(1-\gamma)^2 \varepsilon})$, $T = O(\frac{C_{\text{push}}^2 \iota^2}{(1-\gamma)^{18} \varepsilon^4})$, $N = O(\frac{C_{\text{push}}^2 |\mathcal{A}|^2 \iota^2}{(1-\gamma)^{18} \alpha^2 \varepsilon^6})$,*

$M = m\left(\frac{(1-\gamma)^9 \alpha \epsilon^3}{C_{\text{push}} |\mathcal{A}|}, \frac{\delta}{K T N N_{\text{weight}}}\right)$, $N_{\text{weight}} = O\left(\frac{C_{\text{push}} |\mathcal{A}| \log(|\bar{\mathcal{W}}| \delta^{-1})^\iota}{(1-\gamma)^{11} \epsilon^3}\right)$, where $\iota = \log\left(\frac{C_{\text{push}}}{(1-\gamma)^2 \epsilon}\right)$. Then Algorithm 3 produces policy π , such that with probability at least $1 - \delta$,

$$V^* - V^\pi \leq \frac{C_{\text{push}} \iota}{(1-\gamma)^6 \epsilon} \mathcal{E}(\Pi, \Pi) + \epsilon.$$

547 **Sample Complexity:** If $m(\epsilon, \delta) = \frac{\log |\mathcal{W}|}{\epsilon^2} \log \frac{1}{\delta}$ for some measure of weak learning complexity $|\mathcal{W}|$,
 548 then the algorithm uses at most $K(TNM + N_{\text{weight}}) \leq \tilde{O}\left(\frac{C_{\text{push}}^7 |\mathcal{A}|^4 \log(|\mathcal{W}|)}{(1-\gamma)^{53} \alpha^4 \epsilon^{17}} + \frac{C_{\text{push}}^2 |\mathcal{A}| \log(|\bar{\mathcal{W}}|)}{(1-\gamma)^{13} \epsilon^4}\right)$
 549 episodes.

550 *Proof of Theorem 3.5.* We follow a similar argument to the proof of Theorem 3.2, and obtain:

$$\frac{1}{K} \sum_{k=1}^K V^{\pi^*} - V^{\pi^k} = \underbrace{\frac{1}{K} \sum_{k=1}^K V^{\pi^*} - \bar{V}_{\hat{b}^k}^{\pi^*}}_{I_1} + \underbrace{\frac{1}{K} \sum_{k=1}^K \bar{V}_{\hat{b}^k}^{\pi^*} - \bar{V}_{\hat{b}^k}^{\pi^k}}_{I_2} + \underbrace{\frac{1}{K} \sum_{k=1}^K \bar{V}_{\hat{b}^k}^{\pi^k} - V^{\pi^k}}_{I_3},$$

551 where $\bar{\pi}^*$ be the truncated policy counterpart of π^* . We have $I_1 \leq 0$ follows by Claim 1.

552 For I_3 ,

$$\begin{aligned} \bar{V}_{\hat{b}^k}^{\pi^k} - V^{\pi^k} &= \mathbb{E}^{\pi^k} \left[\frac{2}{1-\gamma} \hat{w}^k(x | x', a) \right] \\ &\leq \frac{6}{1-\gamma} \mathbb{E}^{\pi^k} [w^k(x | x', a)] + \frac{4}{1-\gamma} \mathbb{E}^{\pi^k} \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \right], \end{aligned}$$

553 where the inequality follows from Lemma L.5,

554 Let $g^i = \mathbb{E}^{\pi^i} \left[\left(\sqrt{\hat{w}^i(x | x', a')} - \sqrt{w^i(x | x', a')} \right)^2 \mid a' = \pi^i(x') \right]$, apply Lemma L.6

$$\begin{aligned} &\sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \right] \\ &\leq \sqrt{2C_{\text{push}} \log(K) \sum_{k \in [K]} \sum_{i < k} \mathbb{E}^{\pi^i} \left[\left(\mathbb{E} \left[\left(\sqrt{\hat{w}^i(x | x', a')} - \sqrt{w^i(x | x', a')} \right)^2 \mid a' = \pi^i(x') \right] \right)^2 \right]} + 4C_{\text{push}} \\ &\leq \sqrt{2C_{\text{push}} \log(2K) \sum_{k=1}^K \sum_{i < k} \mathbb{E}^{\pi^i \circ \pi^k} \left[\left(\sqrt{\hat{w}^i(x | x', a')} - \sqrt{w^i(x | x', a')} \right)^4 \right]} + 4C_{\text{push}} \\ &\leq \sqrt{2C_{\text{push}} |\mathcal{A}| \log(2K) \sum_{k=1}^K \sum_{i < k} \mathbb{E}^{\pi^i \circ \pi_{\text{unif}}} \left[\left(\sqrt{\hat{w}^i(x | x', a')} - \sqrt{w^i(x | x', a')} \right)^4 \right]} + 4C_{\text{push}} \\ &\leq \sqrt{8C_{\text{push}} |\mathcal{A}| \log(K) \sum_{k=1}^K \sum_{i < k} \mathbb{E}^{\pi^i \circ \pi_{\text{unif}}} \left[\left(\sqrt{\hat{w}^i(x | x', a')} - \sqrt{w^i(x | x', a')} \right)^2 \right]} + 4C_{\text{push}} \\ &\leq \sqrt{8C_{\text{push}} |\mathcal{A}| K^2 \log(K) \epsilon_{\text{weight}}^2} + 4C_{\text{push}}. \end{aligned} \tag{7}$$

555 Next, we apply the following lemma to bound $\sum_{k=1}^K \mathbb{E}^{\pi^k} [w^k(x | x', a)]$,

556 **Lemma E.1.** Let π^1, \dots, π^K be the set of policy obtained from Algorithm 3, we have that

$$\sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\frac{P(x' | x, a)}{\sum_{i=1}^{k-1} d^{\pi^i}(x') + P(x' | x, a)} \right] \leq 4C_{\text{push}} \log(2K)$$

557 Thus, we have the following bound for I_3 :

$$\frac{1}{K} \sum_{k=1}^K \bar{V}_{\hat{b}^k}^{\pi^k} - V^{\pi^k} \leq \frac{24C_{\text{push}} \log(2K)}{(1-\gamma)K} + \frac{4}{1-\gamma} \sqrt{8C_{\text{push}} |\mathcal{A}| \log(K) \varepsilon_{\text{weight}}^2} + \frac{16C_{\text{push}}}{(1-\gamma)K}$$

558 For I2, similar to Theorem 3.2, we apply Lemma L.10, which requires a guarantee for the gradient
559 domination, we provide the following lemma:

560 **Lemma E.2.** For any $k \in [K]$, let $\kappa = \frac{k-1}{\zeta(1-\gamma)}$, $\tau = \frac{2(k-1)}{\zeta(1-\gamma)^3} \mathcal{E}(\mathbb{I}, \Pi) + \frac{4}{(1-\gamma)^2} \sqrt{\frac{(k-1)|\mathcal{A}| \varepsilon_{\text{weight}}^2}{\zeta}}$ +
561 $\frac{\zeta}{(1-\gamma)^2}$. We have $\bar{V}_{\hat{b}^k}^{\pi^{k-1}}$ is $(\kappa, \tau, \mathbb{I}, \Pi)$ -gradient dominated. That is, for any $\pi \in \mathbb{I}$,

$$\max_{\bar{\pi} \in \bar{\Pi}_k} \bar{V}_{\hat{b}^k}^{\bar{\pi}} - \bar{V}_{\hat{b}^k}^{\pi} \leq \kappa \max_{\pi \in \mathbb{I}} (\nabla \bar{V}_{\hat{b}^k}^{\pi})^\top (\pi' - \pi) + \tau.$$

562 Put everything together, we have

$$\begin{aligned} V_{d_0}^{\pi^*} - V_{d_0}^{\pi^{k'}} &\leq \frac{2K}{(1-\gamma)^4 \zeta} \mathcal{E}(\mathbb{I}, \Pi) + \frac{8K^2}{\zeta^2 (1-\gamma)^6 T} + \frac{4\varepsilon_w K |\mathcal{A}|}{\zeta (1-\gamma)^4 \alpha} + \frac{8K |\mathcal{A}|}{\zeta (1-\gamma)^4 \alpha \sqrt{N}} \\ &\quad + \frac{4}{(1-\gamma)^3} \sqrt{\frac{K |\mathcal{A}| \varepsilon_{\text{weight}}^2}{\zeta}} + \frac{\zeta}{(1-\gamma)^3} + \frac{4}{(1-\gamma)^2} \sqrt{8C_{\text{push}} |\mathcal{A}| \log(2K) \varepsilon_{\text{weight}}^2} \\ &\quad + \frac{48C_{\text{push}} \log(2K)}{(1-\gamma)^2 K}, \end{aligned}$$

563 we conclude the proof by choosing $\zeta = (1-\gamma)^3 \varepsilon$. \square

564 E.1 Support proof for pushforward coverability

565 *Proof of Lemma E.2.* Let $\bar{\pi} \in \bar{\Pi}_k$. By the performance difference lemma (Lemma L.1), we have the
566 following:

$$\begin{aligned} \bar{V}_{\hat{b}^k}^{\bar{\pi}} - \bar{V}_{\hat{b}^k}^{\pi} &= \frac{1}{1-\gamma} \sum_{x,a} \mathbb{E}_{x \sim d^{\bar{\pi}}} \left[\left(\bar{Q}_{\hat{b}^k}^{\bar{\pi}}(x,a) - \bar{V}_{\hat{b}^k}^{\pi}(x) \right) \mathbb{1}\{x \in \mathcal{B}^k\} \right] \\ &\quad + \frac{1}{1-\gamma} \mathbb{E}_{x \sim d^{\bar{\pi}}} \left[\left(\bar{Q}_{\hat{b}^k}^{\bar{\pi}}(x,a) - \bar{V}_{\hat{b}^k}^{\pi}(x) \right) \mathbb{1}\{x \notin \mathcal{B}^k, x \neq \mathfrak{t}\} \right], \end{aligned} \quad (8)$$

567 Due to the construction of the extended MDP, $\bar{Q}_{\hat{b}^k}^{\bar{\pi}}(\mathfrak{t}, a) = 0, \forall a \in \bar{\mathcal{A}}$, which remove term with $x = \mathfrak{t}$.

568 For the first term in Eq. (8):

$$\begin{aligned} &\frac{1}{1-\gamma} \sum_{x,a} d^{\bar{\pi}}(x) \bar{\pi}(a | x) \left(\bar{Q}_{\hat{b}^k}^{\bar{\pi}}(x,a) - \bar{V}_{\hat{b}^k}^{\pi}(x) \right) \mathbb{1}\{x \in \mathcal{B}\} \\ &\leq \frac{k-1}{\zeta(1-\gamma)^2} \sum_{x,a} d^{\pi}(x) \bar{\pi}(a | x) \left(\bar{Q}_{\hat{b}^k}^{\bar{\pi}}(x,a) - \bar{V}_{\hat{b}^k}^{\pi}(x) \right) \mathbb{1}\{x \in \mathcal{B}\} \\ &\leq \frac{k-1}{\zeta(1-\gamma)^2} \sum_x d^{\pi}(x) \left(\max_a \bar{Q}_{\hat{b}^k}^{\bar{\pi}}(x,a) - \bar{V}_{\hat{b}^k}^{\pi}(x) \right) \\ &= \frac{k-1}{\zeta(1-\gamma)^2} \max_{\pi_0} \sum_{x,a} d^{\pi}(x) \pi_0(a | x) \left(\bar{Q}_{\hat{b}^k}^{\bar{\pi}}(x,a) - \bar{V}_{\hat{b}^k}^{\pi}(x) \right) \\ &= \frac{k-1}{\zeta(1-\gamma)} \max_{\pi_0} (\nabla \bar{V}_{\hat{b}^k}^{\bar{\pi}})^\top (\pi_0 - \pi) \\ &\leq \frac{k-1}{\zeta(1-\gamma)} \max_{\pi' \in \mathbb{I}} (\nabla \bar{V}_{\hat{b}^k}^{\bar{\pi}})^\top (\pi' - \pi) + \frac{2(k-1)}{\zeta(1-\gamma)^3} \mathcal{E}(\mathbb{I}, \Pi) \end{aligned}$$

569 where the first inequality follows by $\frac{d^{\bar{\pi}}(x)}{d^{\pi}(x)} \leq \frac{d^{\bar{\pi}}(x)}{d^{\pi^k}(x)(1-\gamma)} \leq \frac{k-1}{\zeta(1-\gamma)}$, the last inequality follows a
570 similar argument in Lemma D.1.

571 To bound second term in Eq. (8), we consider the following:

$$\begin{aligned}
& \mathbb{E}_{x \sim d^\pi} \left[\left(\bar{Q}_{b^k}^\pi(x, \mathbf{t}) - \bar{V}_{\hat{b}^k}^\pi(x) \right) \mathbb{1}\{x \notin \mathcal{B}, x \neq \mathbf{t}\} \right] \\
& \leq \mathbb{E}_{x \sim d^\pi} \left[\left(\bar{Q}_{b^k}^\pi(x, \mathbf{t}) - \frac{2}{1-\gamma} \mathbb{E}_{x', a'}^\pi [\hat{w}^k(\tilde{x} | x', a') | \tilde{x} = x] \right) \mathbb{1}\{x \notin \mathcal{B}, x \neq \mathbf{t}\} \right] \\
& = \mathbb{E}_{x \sim d^\pi} \left[\underbrace{\left(\bar{Q}_{b^k}^\pi(x, \mathbf{t}) - \frac{2}{1-\gamma} \mathbb{E}_{x', a'}^\pi [w^k(\tilde{x} | x', a') | \tilde{x} = x] \right) \mathbb{1}\{x \notin \mathcal{B}, x \neq \mathbf{t}\}}_{A_1} \right] \\
& \quad + \frac{2}{1-\gamma} \underbrace{\mathbb{E}_{x \sim d^\pi} \left[\left(\mathbb{E}_{x', a'}^\pi [w^k(\tilde{x} | x', a') - \hat{w}^k(\tilde{x} | x', a') | \tilde{x} = x] \right) \mathbb{1}\{x \notin \mathcal{B}, x \neq \mathbf{t}\} \right]}_{A_2}
\end{aligned}$$

572 For any $x \notin \mathcal{B}^k$ (i.e., $\frac{1}{d^{p^k}(x)} \geq \frac{k-1}{\zeta}$), any $P(\tilde{x} | x', a') \geq \zeta$, we have

$$b^k(\tilde{x} | x', a') = \frac{2}{1-\gamma} w^k(\tilde{x} | x', a') = \frac{2}{1-\gamma} \frac{P(\tilde{x} | x', a')}{\sum_{i < k} d^{\pi^i}(\tilde{x}) + P(\tilde{x} | x', a')} \geq \frac{1}{1-\gamma},$$

573 also, due to the dynamics of the extended MDP, we have $\bar{Q}_{b^k}^\pi(x, \mathbf{t}) = \frac{1}{1-\gamma}$, therefore, we can bound
574 A_1 as

$$\begin{aligned}
A_1 & \leq \mathbb{E}_{x', a'}^\pi \left[\frac{1}{1-\gamma} \mathbb{1}\{x \notin \mathcal{B}^k, P(\tilde{x} | x', a') < \zeta\} \right] \\
& \leq \frac{1}{1-\gamma} \sum_{x', a' \in \mathcal{X} \times \mathcal{A}} d^\pi(x', a') P(\tilde{x} | x', a') \mathbb{1}\{x \notin \mathcal{B}^k, P(\tilde{x} | x', a') < \zeta\} \\
& < \frac{\zeta}{1-\gamma} \sum_{x', a' \in \mathcal{X} \times \mathcal{A}} d^\pi(x', a') \\
& = \frac{\zeta}{1-\gamma}
\end{aligned} \tag{9}$$

575 By Lemma L.5, we have

$$\mathbb{E}_{x', a'}^\pi [w^k(\tilde{x} | x', a') - \hat{w}^k(\tilde{x} | x', a') | \tilde{x} \neq \mathbf{t}] \leq 2 \sqrt{\mathbb{E}_{x', a'}^\pi \left[\left(\sqrt{w^k(\tilde{x} | x', a')} - \sqrt{\hat{w}^k(\tilde{x} | x', a')} \right)^2 | \tilde{x} \neq \mathbf{t} \right]}$$

576 We complete the proof by bounding A_2 as follows:

$$\begin{aligned}
A_2 &\leq 2\mathbb{E}_{x \sim d^\pi} \sqrt{\mathbb{E}_{x', a'}^{\pi} \left[\left(\sqrt{w^k(x | x', a')} - \sqrt{\widehat{w}^k(x | x', a')} \right)^2 \mid x \neq \mathfrak{t} \right]} \quad (\text{Lemma L.5}) \\
&\leq 2\sqrt{\mathbb{E}_{x \sim d^\pi} \left[\mathbb{E}_{x', a'}^{\pi} \left[\left(\sqrt{w^k(x | x', a')} - \sqrt{\widehat{w}^k(x | x', a')} \right)^2 \mid x \neq \mathfrak{t} \right] \right]} \quad (\text{Jensen's inequality}) \\
&\leq 2\sqrt{\sum_{x \in \mathcal{X}: x \neq \mathfrak{t}} d^\pi(x) \sum_{\substack{x', a' \in \mathcal{X} \times \mathcal{A}: \\ x' \neq \mathfrak{t}, a' \neq \mathfrak{t}}} d^\pi(x') \pi(a' | x') P(x | x', a') \left(\sqrt{w^k(x | x', a')} - \sqrt{\widehat{w}^k(x | x', a')} \right)^2} \\
&\leq 2\sqrt{\frac{k-1}{\zeta} \sum_{\substack{x \in \mathcal{X}: \\ x \neq \mathfrak{t}}} \sum_{\substack{x', a' \in \mathcal{X} \times \mathcal{A}: \\ x' \neq \mathfrak{t}, a' \neq \mathfrak{t}}} d^{p^k}(x') \pi(a' | x') P(x | x', a') \left(\sqrt{w^k(x | x', a')} - \sqrt{\widehat{w}^k(x | x', a')} \right)^2} \\
&\leq 2\sqrt{\frac{(k-1)|\mathcal{A}|}{\zeta} \sum_{\substack{x \in \mathcal{X}: \\ x \neq \mathfrak{t}}} \sum_{\substack{x', a' \in \mathcal{X} \times \mathcal{A}: \\ x' \neq \mathfrak{t}, a' \neq \mathfrak{t}}} d^{p^k}(x') \pi_{\text{unif}}(a' | x') P(x | x', a') \left(\sqrt{w^k(x | x', a')} - \sqrt{\widehat{w}^k(x | x', a')} \right)^2} \\
&= 2\sqrt{\frac{(k-1)|\mathcal{A}|}{\zeta(k-1)} \sum_{i < k} \mathbb{E}^{\pi^i \circ \pi_{\text{unif}}} \left[\left(\sqrt{w^k(x | x', a')} - \sqrt{\widehat{w}^k(x | x', a')} \right)^2 \right]} \\
&= 2\sqrt{\frac{(k-1)|\mathcal{A}| \varepsilon_{\text{weight}}^2}{\zeta}}, \tag{10}
\end{aligned}$$

577 where in the third inequality, since $x \neq \mathfrak{t}$, we have $x' \neq \mathfrak{t}, a' \neq \mathfrak{t}$ for the state x' and action a' in the
578 previous step before reaching x . \square

579 *Proof of Lemma E.1.* Let $\widetilde{d}^k(x) = \sum_{i=1}^{k-1} d^{\pi^i}(x)$, let $\mu \in \Delta(\mathcal{X})$ attain the value of C_{push} , with $\epsilon = 1$,
580 and $\delta = C_{\text{push}}$, applying Lemma E.3:

$$\begin{aligned}
\mathbb{E}^{\pi^k} \left[\frac{P(x' | x, a)}{\widetilde{d}^k(x') + P(x' | x, a)} \right] &\leq \mathbb{E}^{\pi^k} \left[\frac{P(x' | x, a)}{\widetilde{d}^k(x') + C_{\text{push}} \mu(x')} \right] + C_{\text{push}} \cdot \mathbb{E}^{\pi^k} \left[\frac{\mu(x')}{\widetilde{d}^k(x') + C_{\text{push}} \mu(x')} \right] \\
&\leq 2C_{\text{push}} \cdot \mathbb{E}^{\pi^k} \left[\frac{\mu(x')}{\widetilde{d}^k(x') + C_{\text{push}} \mu(x')} \right].
\end{aligned}$$

581 Since $\sup_{\pi \in \Pi} d^\pi(x) \leq \sup_{x' \in \mathcal{X}, a \in \mathcal{A}} P(x | x', a) \leq C_{\text{push}} \mu(x)$ for all $x \in \mathcal{X}$, Lemma L.3 implies
582 that

$$\sum_{x \in \mathcal{X}} \sum_{k=1}^K \mu(x) \frac{d^{\pi^k}(x)}{\widetilde{d}^k(x) + C_{\text{push}} \mu(x)} \leq 2 \log(2K),$$

583 which concludes the proof. \square

Lemma E.3. For any policy π , state weight vectors $d, \mu \in \mathbb{R}_+^{\mathcal{X}}$, and constants $\varepsilon, \delta > 0$, we have:

$$\mathbb{E}^\pi \left[\frac{P(x' | x, a)}{d(x') + \varepsilon \cdot P(x' | x, a)} \right] \leq \mathbb{E}^\pi \left[\frac{P(x' | x, a)}{d(x') + \delta \cdot \mu(x')} \right] + \frac{\delta}{\varepsilon} \cdot \mathbb{E}^\pi \left[\frac{\mu(x')}{d(x') + \delta \cdot \mu(x')} \right]$$

Proof.

$$\begin{aligned}
& \mathbb{E}^\pi \left[\frac{P(x' | x, a)}{d(x') + \varepsilon \cdot P(x' | x, a)} \right] - \mathbb{E}^\pi \left[\frac{P(x' | x, a)}{d(x') + \delta \cdot \mu(x')} \right] \\
&= \sum_{x, a, x'} d^\pi(x, a) \frac{P(x' | x, a) (\delta \cdot \mu(x') - \varepsilon \cdot P(x' | x, a))}{(d(x') + \varepsilon \cdot P(x' | x, a)) (d(x') + \delta \cdot \mu(x'))} \\
&\leq \sum_{x, a, x'} d^\pi(x, a) \frac{P(x' | x, a) (\delta \cdot \mu(x'))}{(d(x') + \varepsilon \cdot P(x' | x, a)) (d(x') + \delta \cdot \mu(x'))} \\
&\leq \frac{\delta}{\varepsilon} \sum_{x, a, x'} d^\pi(x, a) \frac{P(x' | x, a) \mu(x')}{d(x') + \delta \cdot \mu(x')} \\
&= \frac{\delta}{\varepsilon} \cdot \mathbb{E}^\pi \left[\frac{\mu(x')}{d(x') + \delta \cdot \mu(x')} \right].
\end{aligned}$$

584

□

585 **Weight estimation.** To estimate the weight function, consider the following setting. Let \mathcal{Z} be
586 a set. We receive samples $z_\mu^1, \dots, z_\mu^n \in \mathcal{Z}$ and $z_\nu^1, \dots, z_\nu^n \in \mathcal{Z}$, where $z_\mu^t \sim \mu^t \in \Delta(\mathcal{Z})$ and
587 $z_\nu^t \sim \nu^t \in \Delta(\mathcal{Z})$. The distributions μ^t and ν^t can be chosen in an adaptive fashion based on
588 $z_\mu^1, z_\nu^1, \dots, z_\mu^{t-1}, z_\nu^{t-1}$. We define $\mu = \frac{1}{n} \sum_{t=1}^n \mu^t$ and $\nu = \frac{1}{n} \sum_{t=1}^n \nu^t$, and our goal is to estimate
589 the density ratio

$$w^*(z) := \frac{\mu(z)}{\nu(z)}.$$

590 We assume that $\|w_\star\|_\infty \leq B$, and assume access to a realizable weight function class \mathcal{W} with
591 $w^\star \in \mathcal{W}$. We consider the estimator

$$\hat{w} := \arg \max_{w \in \mathcal{W}} \widehat{\mathbb{E}}_\mu[\log(w)] - \widehat{\mathbb{E}}_\nu[w], \quad (11)$$

592 where $\widehat{\mathbb{E}}_\mu[\cdot]$ denotes the empirical expectation with respect to z_μ^1, \dots, z_μ^n and $\widehat{\mathbb{E}}_\nu[\cdot]$ denotes the
593 empirical expectation with respect to z_ν^1, \dots, z_ν^n . The following theorem gives a finite-sample bound
594 for this estimator, which may be of independent interest.

595 **Lemma E.4.** *Suppose that $w^\star \in \mathcal{W}$ and $\sup_{w \in \mathcal{W}} \|w\|_\infty \leq B$. The estimator in Eq. (11) ensures*
596 *that with probability at least $1 - \delta$,*

$$D_{\mathbb{H}, \nu}^2(\hat{w}, w^\star) \leq \frac{20B \log(|\mathcal{W}| \delta^{-1})}{n},$$

597 where $D_{\mathbb{H}, \nu}^2(w, w') := \mathbb{E}_\nu \left[\left(\sqrt{w} - \sqrt{w'} \right)^2 \right]$.

598 To state the guarantee for the weight estimation, we use $\pi \circ \pi'$ to denote the policy that plays the
599 policy π' one step after π is “terminated” Specifically, under the episodic reset interpretation of a
600 discounted MDP, we execute policy π and terminate with probability $1 - \gamma$; upon termination, we
601 then execute π' one additional step. For example, in the following weight estimate lemma, we use
602 $\pi^i \circ \pi_{\text{unif}}$, meaning that we run π^i with probability γ , and with probability $1 - \gamma$, we run one time
603 step with π_{unif} and terminate.

604 **Lemma E.5.** *For any $k \in [K]$, for any $\varepsilon_{\text{weight}}, \delta \in (0, 1)$, distribution $p^k \in \Delta(\mathbb{I})$ and $\pi^1, \dots, \pi^1 \in$
605 \mathbb{I} , Algorithm 4 ensures that with probability at least $1 - \delta$, the output \hat{w}^k satisfies*

$$\frac{1}{k-1} \sum_{i < k} \mathbb{E}^{\pi^i \circ \pi_{\text{unif}}} \left[\left(\sqrt{\widehat{w}^i(x | x', a')} - \sqrt{w^i(x | x', a')} \right)^2 \right] \leq \varepsilon_{\text{weight}}^2,$$

606 and does so using at most $N_{\text{weight}} = 40 \frac{\log(|\overline{\mathcal{W}}| \delta^{-1})}{\varepsilon_{\text{weight}}^2}$ episodes.

607 *Proof.* This proof follows similarly to Lemma J.7 of Amortila et al. (2024). Let $\bar{w}^k := k \cdot w^k$ and let
 608 $\check{w}^k := k \cdot \widehat{w}^k$. Observe that solving the optimization problem in Line 10 of Algorithm 5 is equivalent
 609 to solving the optimization problem in over the class $t \cdot \mathcal{W}_h$, which has $\|w'\|_{\infty} \leq t$ for all $w' \in \mathcal{W}$.
 610 As such, we can appeal to Lemma E.5 with

$$\mu(x' | x, a) = P(x' | x, a), \quad \nu(x' | x, a) = \frac{1}{k} \left(\sum_{i < k} d_{\pi^i}(x') + P(x' | x, a) \right),$$

611 and

$$\omega(x, a) = \frac{1}{2} \left(d^{p^k}(x, a) + \frac{1}{k-1} \sum_{i < k} d^{\pi^i \circ \pi_{\text{unif}}}(x, a) \right).$$

612 Under the weight assumption, we have

$$\frac{\mu(x' | x, a)}{\nu(x' | x, a)} = \bar{w}(x' | x, a) \in k \cdot \overline{\mathcal{W}},$$

613 so Lemma E.5 imply that

$$\mathbb{E}_{(x,a) \sim \omega} \left[\left(\sqrt{\check{w}_h^k(x' | x, a)} - \sqrt{\bar{w}^k(x' | x, a)} \right)^2 \right] \leq \frac{20k \log(|\overline{\mathcal{W}}| \delta^{-1})}{n},$$

614 or equivalently,

$$\mathbb{E}_{(x,a) \sim \omega} \left[\left(\sqrt{\widehat{w}^k(x' | x, a)} - \sqrt{w^k(x' | x, a)} \right)^2 \right] \leq \frac{20 \log(|\overline{\mathcal{W}}| \delta^{-1})}{n}.$$

615

□

616 **F Reward-free Exploration**

617 In this section, we provide a result on a reward-free setting. We note that the Algorithm 3 can be
 618 directly used in a reward-free setting by simply setting $r(x, a) = 0, \forall x \in \mathcal{X}, a \in \mathcal{A}$. Which results
 619 in an algorithm that only optimizes to maximize the total bonus. We can safely ignore the factor $\frac{1}{1-\gamma}$
 620 since this is a constant. In other word, the optimization at each iteration k is

$$\pi = \operatorname{argmax}_{\pi} \mathbb{E}^{\pi} [w^k(x' | x, a)].$$

621 So basically, the policy operates in the same MDP with the reward $r = w^k$ in iteration k instead.

The policy cover obtained from Algorithm 3 can be later used on a downstream task with a specific-reward function. To see this, we first define the L_1 -coverage from Amortila et al. (2024) as follows:

$$\Psi_{\varepsilon}^M(p) = \sup_{\pi \in \Pi} \mathbb{E}^{\pi} \left[\frac{d^{\pi}(x, a)}{d^p(x, a) + \varepsilon \cdot d^{\pi}(x, a)} \right].$$

622 It can be shown that L_1 coverage enables downstream task optimization by using a standard offline
 623 optimizer, where the sample complexity scales with the $\Psi_{\varepsilon}(p)$ (Amortila et al., 2024). Furthermore,
 624 we can relate $\Psi_{\varepsilon}(p)$ with the pushforward coverability by

$$\Psi_{\varepsilon}^M(p) \leq |\mathcal{A}| \cdot \Psi_{\text{push};\varepsilon}^M(p),$$

625 where we define $\Psi_{\text{push};\varepsilon}^M(p) = \sup_{\pi \in \Pi} \mathbb{E}^{\pi} \left[\frac{P^M(x|x', a')}{d^p(x) + \varepsilon \cdot P(x|x', a')} \right]$.

626 We show that $\Psi_{\varepsilon}(p)$ is bounded and depends on C_{push} , where p is the policy cover obtained from
 627 Algorithm 3.

628 **Theorem F.1.** *In the reward-free setting, suppose Algorithm 3 is instantiated with the following*
 629 *parameters: $T = O(\frac{1}{\varepsilon^3(1-\gamma)^4})$, $N = O(\frac{|\mathcal{A}|^2}{(1-\gamma)^4 \alpha^2 \varepsilon^4})$, $M = m \left(\frac{(1-\gamma)^2 \alpha \varepsilon^2}{|\mathcal{A}|}, \frac{\delta}{KNT} \right)$, $N_{\text{weight}} =$*
 630 *$O(\frac{C_{\text{push}} |\mathcal{A}| \log(|\mathcal{W}| \delta^{-1})}{(1-\gamma)^2 \varepsilon^2})$, we have with probability $1 - \delta$,*

$$\Psi_{\text{push};h,\varepsilon}(p_h) \leq \frac{9\mathcal{E}(\Pi, \Pi)}{(1-\gamma)^2 \varepsilon^2} + \frac{173}{1-\gamma} C_{\text{push}} \log(\varepsilon^{-1}).$$

631 *Proof of Theorem F.1.* For the analysis of the reward-free setting, we use the following truncated
 632 policy class $\bar{\Pi}_K$ such that for any policy $\pi \in \bar{\Pi}_K$:

$$\bar{\pi}(x) = \begin{cases} \pi(x), & \frac{\bar{d}^{\pi}(x)}{\bar{d}^{p^k}(x)} \leq K \\ \mathbf{t}, & \text{otherwise,} \end{cases}$$

633 We construct the *extended* MDP such that when a policy pick the terminal action \mathbf{t} it deterministically
 634 transition to the terminal state \mathbf{t} and immediately receive zero reward. Let r be any reward function
 635 and the reward function at each time step is bounded by 1, so whenever the truncated policy pick \mathbf{t} , it
 636 transitions to the self-absorb state, therefore, we have the following:

$$\sup_{\pi \in \Pi} \mathbb{E}^{\pi} [r] - \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^{\pi} [r] \leq \frac{1}{1-\gamma} \sup_{\pi \in \bar{\Pi}_k} \bar{P}^{\pi} \left[\frac{\bar{d}^{\pi}(x)}{\bar{d}^{p^k}(x)} > K, x \neq \mathbf{t} \right]. \quad (12)$$

$$\begin{aligned}
\sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\sum_{\ell=1}^{h-1} r_\ell \right] &= \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\mathbb{E} \left[\frac{P(x | x', a')}{\sum_{i < k} d^{\pi^i}(x) + P(x | x', a')} \middle| x', a' \right] \mathbb{1}\{x', a' \neq \mathfrak{t}\} \right] \\
&= \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\bar{\mathbb{E}} \left[\frac{\bar{P}(x | x', a')}{\sum_{i < k} \bar{d}^{\pi^i}(x) + \bar{P}(x | x', a')} \middle| x', a' \right] \mathbb{1}\{x', a' \neq \mathfrak{t}\} \right] \quad (13) \\
&= \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x', a')}{\sum_{i < k} d^{\pi^i}(x) + \bar{P}(x | x', a')} \mathbb{1}\{x', a' \neq \mathfrak{t}\} \right] \\
&= \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x', a')}{\sum_{i < k} d^{\pi^i}(x) + \bar{P}(x | x', a')} \mathbb{1}\{x' \neq \mathfrak{t}\} \right]
\end{aligned}$$

637 where we use $\varepsilon = \frac{1}{k-1}$ in the last line. Furthermore, we have the following result in the discounted
638 MDP

Lemma F.2. *For all $K > 0$, it holds that ¹*

$$\sup_{\pi \in \bar{\Pi}_K} \bar{P}^\pi \left[\frac{\bar{d}^\pi(x)}{\bar{d}^{\pi^k}(x)} > K, x \neq \mathfrak{t} \right] \leq 2 \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x, a)}{\sum_{i < K+1} \bar{d}^{\pi^i}(x) + \bar{P}(x | x, a)} \mathbb{1}\{x \neq \mathfrak{t}\} \right]$$

639 From Eqs. (12) and (13), Lemma F.2, we have that:

$$\sup_{\pi \in \bar{\Pi}} \bar{\mathbb{E}}^\pi [w^K] \leq \frac{3}{1-\gamma} \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x, a)}{\sum_{i < K+1} \bar{d}^{\pi^i}(x) + \bar{P}(x | x, a)} \mathbb{1}\{x \neq \mathfrak{t}\} \right] \quad (14)$$

640 and hence it suffices to bound the quantity on the right-hand side. Next, note that for all $t \in [T]$, we
641 have that

$$\sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x', a')}{\bar{d}^k(x) + \bar{P}(x | x', a')} \mathbb{1}\{x \neq \mathfrak{t}\} \right] \leq \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x', a')}{\bar{d}^{k-1}(x) + \bar{P}(x | x', a')} \mathbb{1}\{x \neq \mathfrak{t}\} \right]$$

642 Therefore,

$$K \cdot \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x', a')}{\bar{d}^K(x) + \bar{P}(x | x', a')} \mathbb{1}\{x \neq \mathfrak{t}\} \right] \leq \sum_{k=1}^K \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x', a')}{\bar{d}^k(x) + \bar{P}(x | x', a')} \mathbb{1}\{x \neq \mathfrak{t}\} \right] \quad (15)$$

643 Now, note that if $x \neq \mathfrak{t}$, the dynamics of $\bar{\mathcal{M}}$ imply that we must have $x', a' \neq \mathfrak{t}$ as well. In this case,
644 we have

$$\frac{\bar{P}(x | x', a')}{\bar{d}_h^t(x) + \bar{P}(x | x', a')} = w^k(x | x', a'),$$

645 since $\bar{P}(\cdot | x', a') = P(\cdot | x', a')$ with $x', a' \neq \mathfrak{t}$, and since $\bar{d}_h^{\pi^i}(x) = \bar{d}_h^{\pi^i}(x)$ when $x \neq \mathfrak{t}$ (as the
646 policies π^1, \dots, π^k never take the terminal action). As a result, using Lemma J.2, we have that

¹This is Lemma J.12 of Amortila et al. (2024), although their results apply for episodic MDP, but this also holds for discounted MDP.

$$\begin{aligned}
& \sum_{k=1}^K \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\frac{\bar{P}(x | x', a')}{\bar{d}^k(x) + \bar{P}(x | x', a')} \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&= \sum_{k=1}^K \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[w^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&\leq 3 \sum_{k=1}^K \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\hat{w}^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&\quad + 2 \sum_{k=1}^K \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \mathbb{1}\{x \neq \mathfrak{t}\} \right]
\end{aligned} \tag{16}$$

647 In addition, we have

$$\begin{aligned}
& \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&\leq \sup_{\pi \in \bar{\Pi}_k} \sum_{x \in \mathcal{X}} \sum_{x', a' \in \mathcal{X} \times \mathcal{A}} \bar{d}^\pi(x') \pi(a' | x') P(x | x', a') \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&\leq K \sup_{\pi \in \bar{\Pi}_k} \sum_{x \in \mathcal{X}} \sum_{\substack{x', a' \in \mathcal{X} \times \mathcal{A}: \\ x' \neq \mathfrak{t}, a' \neq \mathfrak{t}}} \bar{d}^{p^k}(x') \pi(a' | x') P(x | x', a') \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&\leq K |\mathcal{A}| \sum_{x \in \mathcal{X}} \sum_{\substack{x', a' \in \mathcal{X} \times \mathcal{A}: \\ x' \neq \mathfrak{t}, a' \neq \mathfrak{t}}} \bar{d}^{\pi^k}(x') \pi_{\text{unif}}(a' | x') P(x | x', a') \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&\leq K |\mathcal{A}| \sum_{x \in \mathcal{X}} \sum_{x', a' \in \mathcal{X} \times \mathcal{A}} \bar{d}^{p^k}(x') \pi_{\text{unif}}(a' | x') P(x | x', a') \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&\leq \frac{K |\mathcal{A}|}{k-1} \sum_{i < k} \bar{\mathbb{E}}^{\pi^i \circ \pi_{\text{unif}}} \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \right] \\
&\leq K |\mathcal{A}| \varepsilon_{\text{weight}}^2
\end{aligned} \tag{17}$$

648 From Eqs. (15) to (17), we have

$$\sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[w^K \right] \leq \frac{3}{K} \sum_{k=1}^K \sup_{\pi \in \bar{\Pi}_K} \bar{\mathbb{E}}^\pi \left[\hat{w}^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] + |\mathcal{A}| \varepsilon_{\text{weight}}^2 \tag{18}$$

649 Next, we proceed to bound $\sum_{k=1}^K \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\hat{w}^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right]$, fix any $k \in [K]$, we
650 have:

$$\begin{aligned}
& \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\hat{w}^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] - \bar{\mathbb{E}}^{\pi^k} \left[\hat{w}^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\
&= \sup_{\pi \in \bar{\Pi}_k} \bar{\mathbb{E}}^\pi \left[\hat{w}^k(x | x', a') \mathbb{1}\{x', a' \neq \mathfrak{t}\} \right] - \bar{\mathbb{E}}^{\pi^k} \left[\hat{w}^k(x | x', a') \mathbb{1}\{x', a' \neq \mathfrak{t}\} \right] \\
&= \bar{V}^{\bar{\pi}} - \bar{V}^{\pi^k} = \varepsilon_{\text{opt}},
\end{aligned}$$

651 where $\bar{V}^{\bar{\pi}} = \max_{\pi' \in \bar{\Pi}_K} \bar{V}^{\pi'}$, now we can apply boosting procedure in [Bruckhim et al. \(2022\)](#) to
652 bound the policy optimization. We first start by showing the following Gradient domination:

653 **Lemma F.3** (Gradient domination). *For any $\bar{\pi} \in \bar{\Pi}$, we have*

$$\bar{V}^{\bar{\pi}} - \bar{V}^{\pi} \leq K \max_{\pi' \in \bar{\Pi}} (\nabla \bar{V}^{\pi'})^\top (\pi' - \pi) + \frac{K}{1-\gamma} \mathcal{E}(\bar{\Pi}, \Pi)$$

654 *Proof.* Apply the performance difference lemma, we have

$$\begin{aligned} \bar{V}^{\bar{\pi}} - \bar{V}^{\pi} &= \sum_{x \in \mathcal{X}} \bar{d}^{\bar{\pi}}(x) \left(\bar{Q}^{\bar{\pi}}(x, \bar{\pi}(x)) - \bar{Q}^{\pi}(x, \pi(x)) \right) \\ &\leq \sum_{x \in \mathcal{X}} \bar{d}^{\bar{\pi}}(x) \left(\bar{Q}^{\bar{\pi}}(x, \bar{\pi}(x)) - \bar{Q}^{\pi}(x, \pi(x)) \right) \mathbb{1}\{\bar{d}^{\bar{\pi}}(x)/\bar{d}^{\pi^k}(x) \leq K\} \\ &\leq \sum_{x \in \mathcal{X}} \bar{d}^{\bar{\pi}}(x) \left(\bar{Q}^{\bar{\pi}}(x, \bar{\pi}(x)) - \bar{Q}^{\pi}(x, \pi(x)) \right) \mathbb{1}\{\bar{d}^{\bar{\pi}}(x)/\bar{d}^{\pi^k}(x) \leq K\} \\ &\leq \sum_{x \in \mathcal{X}} \frac{K}{1-\gamma} \bar{d}^{\bar{\pi}}(x) \left(\bar{Q}^{\bar{\pi}}(x, \bar{\pi}(x)) - \bar{Q}^{\pi}(x, \pi(x)) \right) \mathbb{1}\{\bar{d}^{\bar{\pi}}(x)/\bar{d}^{\pi^k}(x) \leq K\} \\ &\leq K \max_{\pi' \in \bar{\Pi}} (\nabla \bar{V}^{\pi'})^\top (\pi' - \pi) + \frac{K}{1-\gamma} \mathcal{E}(\bar{\Pi}, \Pi) \end{aligned}$$

655 where we use the fact that $\bar{Q}^{\bar{\pi}}(x, \bar{\pi}(x)) - \bar{Q}^{\pi}(x, \pi(x)) \leq 0$ due to the dynamics of the extended MDP. \square

656 For the Internal Boost guarantee, we have that $\varepsilon_0 = \frac{2|A|}{(1-\gamma)\alpha}(\varepsilon_w + \frac{1}{\sqrt{N}})$. Note that the Oracle
657 guarantee run on the true MDP, but we can still apply the result extended MDP, due the the fact that
658 all the policy π we operate in the boosting procedure never pick \mathfrak{t} , therefore, $\bar{Q}^{\bar{\pi}} \equiv Q^{\pi}$. Now, we can
659 apply the Lemma L.10, we have that

$$\sum_{k=1}^K \sup_{\pi \in \bar{\Pi}_k} \mathbb{E}^{\pi} \left[\hat{w}^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] \leq \sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\hat{w}^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] + K\varepsilon_{\text{opt}}, \quad (19)$$

660 where $\varepsilon_{\text{opt}} = \frac{2K^2}{(1-\gamma)^3 T} + \frac{K}{1-\gamma} \mathcal{E}(\bar{\Pi}, \Pi) + \frac{2K|A|}{(1-\gamma)}(\varepsilon_w + \frac{1}{\sqrt{N}})$. For the summation above, we can apply
661 Lemma L.5:

$$\begin{aligned} &\sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\hat{w}^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\ &\leq 3 \sum_{k=1}^K \mathbb{E}^{\pi^k} \left[w^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] + 2 \sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \mathbb{1}\{x \neq \mathfrak{t}\} \right] \end{aligned} \quad (20)$$

662 Finally, note that since $\pi^k \in \Pi$ never select the terminal action (and in particular never reach the
663 terminal state), we have

$$\begin{aligned} \sum_{k=1}^K \mathbb{E}^{\pi^k} \left[w^k(x | x', a') \mathbb{1}\{x \neq \mathfrak{t}\} \right] &= \sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\frac{P(x | x', a')}{\tilde{d}_h^{\mathfrak{t}}(x) + P(x | x', a')} \mathbb{1}\{x \neq \mathfrak{t}\} \right] \\ &= \sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\frac{P(x | x', a')}{\tilde{d}_h^{\mathfrak{t}}(x) + P(x | x', a')} \right] \\ &\leq 4C_{\text{push}} \log(2T), \quad (\text{Lemma E.1}) \end{aligned} \quad (21)$$

664 where in the second inequality, we use the fact that $\bar{\mathbb{E}}^{\pi^k} \equiv \mathbb{E}^{\pi}$ since π^k never pick the terminal action.

665 The weight estimation error follows similar to Eq. (7) in the proof of Theorem 3.5.

$$\sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\left(\sqrt{\hat{w}^k(x | x', a')} - \sqrt{w^k(x | x', a')} \right)^2 \mathbb{1}\{x \neq \mathfrak{t}\} \right] \leq \sqrt{8C_{\text{push}}|A|K^2 \log(K)\varepsilon_{\text{weight}}^2} + 4C_{\text{push}}. \quad (22)$$

666 From Eqs. (18) to (22), and for $\varepsilon = \frac{1}{K}$ we have:

$$\begin{aligned}
\Psi_{\text{push};\varepsilon} &\leq K \cdot \sup_{\pi \in \Pi} \mathbb{E}^\pi \left[\frac{P(x | x', a')}{\sum_{i=1}^K d^{\pi^i}(x) + P(x | x', a')} \right] \\
&\leq \frac{18K^3}{(1-\gamma)^4 T} + \frac{9K^2 \mathcal{E}(\Pi, \Pi)}{(1-\gamma)^2} + \frac{6K^2 |\mathcal{A}|}{(1-\gamma)^2 \alpha} \left(\varepsilon_w + \frac{1}{\sqrt{N}} \right) + \frac{98}{(1-\gamma)} C_{\text{push}} \log(2K) \\
&\quad + \frac{24}{1-\gamma} \sqrt{K^2 C_{\text{push}} |\mathcal{A}| \log(K) \varepsilon_{\text{weight}}^2} + \frac{24 C_{\text{push}}}{1-\gamma} + \frac{2K |\mathcal{A}|}{1-\gamma} \varepsilon_{\text{weight}}^2,
\end{aligned}$$

667 which also complete the proof. □

668 **G RL Boosting in Episodic MDP**

669 **Markov Decision Processes (MDPs).** We consider the episodic MDPs defined as $M =$
 670 $\{H, \mathcal{X}, \mathcal{A}, \mu, \{P_h\}_{h=1}^H\}$, where H is the episode horizon, \mathcal{X} is the state space, \mathcal{A} is the action
 671 space, μ is the initial distribution and $P_h : \mathcal{X} \times \mathcal{A} \rightarrow \Delta(\mathcal{X})$ is the transition kernel at step h . A (pos-
 672 sibly randomized) non-stationary policy is given by $\pi = (\pi_1, \dots, \pi_H)$, where each $\pi_h : \mathcal{X} \rightarrow \Delta(\mathcal{A})$
 673 specifies a distribution over actions. Let Π^H denote the class of such policies constructed from a base
 674 policy class Π , i.e., $\pi_h \in \Pi$ for all $h \in [H]$. At the beginning of each episode, an initial state $x_1 \sim d_0$
 675 is drawn. Then, for each step $h = 1, \dots, H$, the agent selects an action $a_h \sim \pi_h(x_h)$, transitions to
 676 the next state $x_{h+1} \sim P_h(\cdot | x_h, a_h)$, and receives a reward $r_h(x_h, a_h)$.

677 For a policy π , a state x , and $h \in \{1, \dots, H\}$, we define the value function $V_h^\pi : \mathcal{X} \rightarrow \mathbb{R}$ as

$$V_h^\pi(x) = \mathbb{E}^\pi \left[\sum_{h'=h}^H r_{h'}(x_{h'}, a_{h'}) \mid \pi, x_{h'} = x \right],$$

and state-action value function $Q_h^\pi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ as

$$Q_h^\pi(x, a) = \mathbb{E}^\pi \left[\sum_{h'=h}^H r_{h'}(x_{h'}, a_{h'}) \mid \pi, x_{h'} = x, a_{h'} = a \right],$$

678 Furthermore, we define the optimal policy as $\pi^* = \operatorname{argmax}_\pi V_1^\pi(d_0)$, where $V_1^\pi(\mu) =$
 679 $\mathbb{E}_{x_1 \sim d_0}[V_1^\pi(x_1)]$. For simplicity, we denote $V^\pi = V_1^\pi(\mu)$ and $V^{\pi^*} = V_1^{\pi^*}(\mu)$ when it is clear
 680 from the context.

681 We define the *occupancy measure* for layer h as

$$d_h^\pi(x, a) = P^\pi(x_h = x, a_h = a)$$

682 and $d_h^\pi(x) = P^\pi(x_h = x)$. We assume that the cumulative reward is bounded as $\mathbb{E} \left[\sum_{h=1}^H r_h \right] \leq$
 683 R_{\max} . In the episodic MDP, it is common to normalize rewards such that $R_{\max} = 1$ for (Amortila
 684 et al., 2024; Misra et al., 2020). Although we consider an episodic MDP, our results remain
 685 comparable to those in discounted MDP (Bruckhim et al., 2022) by setting the effective horizon
 686 $H \approx \frac{1}{1-\gamma}$, and R_{\max} is bounded by $\frac{1}{1-\gamma}$ where γ is the discount factor. Therefore, we purposely
 687 keep R_{\max} general in our analysis to facilitate comparison across different MDP formulations.

688 Similar to the result in discounted MDPs, our algorithm can work on large-scale MDPs where the
 689 state space \mathcal{X} can be very large, or even infinite. When the environment admits some underlying
 690 structure that allows us to design algorithms that can be computationally and statistically efficient.
 691 We provide a brief description of these structures in episodic MDPs in the following paragraphs.

692 **Block MDPs.** A block MDP models an environment where a large observation space \mathcal{X} can be
 693 summarized into a finite state \mathcal{S} . Specifically, block MDP assumes there exists an emission function
 694 $q : \mathcal{X} \rightarrow \Delta(\mathcal{S})$ that maps the latent states to a distribution over observations. The agent interacts with
 695 the environment by repeatedly generating H -step trajectories $(s_1, x_1, a_1, r_1, \dots, s_H, x_H, a_H, r_H)$,
 696 where $s_1 \sim \mu$, $s_{h+1} \sim P(\cdot | s_h, a_h)$, $x_h \sim q(s_h)$, for all $h \in [H]$. The agent can observe x_h , but the
 697 states s_1, \dots, s_H are hidden from the agent. Tabular MDP is a special case of block MDP where the
 698 observation space coincides with the state space, i.e., $\mathcal{X} = \mathcal{S}$.

699 **Low-rank MDPs.** A low-rank MDP with dimension d is an MDP where the transition kernel
 700 admits a low-rank factorization: $\forall h \in [H]$, there exist $\phi_h : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d$ and $\mu_h : \mathcal{X} \rightarrow \mathbb{R}^d$ such
 701 that $\forall x_h, x_{h+1} \in \mathcal{X}, a_h \in \mathcal{A} : P_h(x_{h+1} | x_h, a_h) = \langle \phi_h(x_h, a_h), \mu_h(x_{h+1}) \rangle$. This structure
 702 enables efficient representation and learning in environments with large or continuous state spaces.

703 Next, we restate several key definitions from the boosting-based reinforcement learning framework
 704 of Bruckhim et al. (2022) with some adjustment in the context of episodic MDPs.

705 **Definition G.1** (Policy Tree). A Policy Tree $\Pi \subseteq \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$ is linear combination of T policy in
 706 $\tilde{\Pi} := \Pi^H \subseteq \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$, for some base policy class Π . Policy at each layer of $\tilde{\pi} \in \tilde{\Pi}$ is a projected
 707 policy $\Gamma(\tilde{\pi})$, where each $\tilde{\pi}$ is a linear combination of N policy in Π . $\Pi \subseteq \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$.

708 We note that in episodic MDP, policy structure is different where the policy is non-stationary which
 709 the policy can be different at different time step. A function policy can be view a concatenation of H
 710 policy for each layer, where each policy belong to some base policy class.

711 **Definition G.2** (Policy completeness). For any $h \in [H]$, distribution μ_h , and any two policy class
 712 \mathbb{I}, \mathbb{II} , we define:

$$\mathcal{E}(\mathbb{I}, \mathbb{II}) = \max_{\pi \in \mathbb{I}} \min_{\pi'_h \in \mathbb{II}} \mathbb{E}_{x_h \sim \mu_h} \left[\max_a Q_h^\pi(x_h, a) - \mathbb{E}_{a \sim \pi'_h} Q_h^\pi(x_h, a) \right]$$

713 The policy completeness measures the ability of the policy class \mathbb{II} at each layer to approximate the
 714 greedy policy with respect to any Q -function induced by a policy in \mathbb{I} .

715 **Definition G.3** (Distribution mismatch). Let $\pi^* = \arg \max_{\pi} V^\pi$. We define the following distribution
 716 mismatch coefficients: $C_\infty = \max_{h \in [H]} \sup_{\pi \in \mathbb{I}} \left\| \frac{d_h^{\pi^*}}{d_h^\pi} \right\|_\infty$.

717 **Further notation.** We denote sequences using the “:” operator, e.g., $\pi_{1:H} = (\pi_1, \dots, \pi_H)$. We use
 718 π_{unif} to denote the uniform policy. We define $\pi \circ_h \pi'$ as the policy that follows π for layers $h' < h$
 719 and follows π' for $h' \geq h$.

720 H Warming up: Boosting in Episodic MDPs

721 The algorithm presented in this section serves as a building block for the policy cover in Appendix I.
 722 These results are a straightforward extension of [Brukhim et al. \(2022\)](#) where we adapt their boosting
 723 approach to episodic MDPs. We note that the algorithm for the policy only needs the first level
 724 boosting algorithm (PSDP Boosting Algorithm 6) as the policy optimization procedure. Nevertheless,
 725 we still provide a similar two-level boosting as in [Brukhim et al. \(2022\)](#), which we provide the pseudo-
 726 code in Algorithm 5. Similar to the boosting procedure in [Brukhim et al. \(2022\)](#), the algorithm runs
 727 for T iterations; at each iteration, it calls a subroutine in Algorithm 6 to obtain a new policy π'_t , which
 728 is then aggregated with the previous policy π_{t-1} to form the new policy π_t . The aggregation is done
 729 via a convex combination with step size $\eta_{1,t}$. Here, we use π_{t-1} to collect samples in Algorithm 6,
 730 and we denote $d_h^{\pi_{t-1}}$ as the state distribution at layer h induced by π_{t-1} . The algorithm outputs the
 731 final policy π^T after T iterations.

Algorithm 5 RL Boosting

- 1: **input:** number of iterations T , number of weak learner N , number of sample episodes M .
 - 2: Initialize policy π^0 arbitrarily.
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Run Algorithm 6 with $\{d_h^{\pi_{t-1}}\}_{h=1}^H$, reward r , number of weak learner N , number of episodes M to obtain π'_t .
 - 5: $\pi_t = (1 - \eta_{1,t})\pi_{t-1} + \eta_{1,t}\pi'_t$.
 - 6: **end for**
 - 7: **return** π_T .
-

Algorithm 6 PSDP Boosting

- 1: **input:** Exploration distribution $\{\rho\}_{h=1}^H$, reward function r , number of weak learner N , number of episodes M
 - 2: **for** $h = H, H - 1, \dots, 1$ **do**
 - 3: Run Algorithm 7 with distribution ρ_h , policy $\pi_{h+1:H}$, number of iterations N , number of episodes M to obtain π_h .
 - 4: **end for**
 - 5: **return** $\pi = (\pi_1, \pi_2, \dots, \pi_H)$.
-

732 Each optimization inside the for loop of PSDP Boosting algorithm, i.e., Algorithm 7, share a similar
 733 flavor with the Internal Boost in [Brukhim et al. \(2022\)](#), where it specifically solves a contextual

734 optimization problem. PSDP Boosting builds the policy from the last layer H to the first layer. At
 735 each layer h , it calls a subroutine Layer Boosting (Algorithm 7) to obtain a policy π_h for that layer.
 736 Layer Boosting runs for N iterations. At each iteration, it collects M episodes by executing $\pi_{h+1:H}$
 737 with the exploration distribution ρ_h . The collected dataset is then modified to form a new dataset
 738 with linear losses (Line 5), which is then fed into the weak learner to obtain a new policy. The new
 739 policy is then aggregated with the previous policy $\tilde{\pi}_{n-1}$ to form the new policy $\tilde{\pi}_n$. The final output
 740 of Layer Boosting is then projected to be a valid policy via the projection operator $\Gamma(\cdot)$ to form the
 741 policy π_h for layer h . We provide the guarantee for the boosting procedure in Algorithm 7 as follow,

742 **Lemma H.1.** Let $\beta = \sqrt{\frac{1}{\alpha N}}$ and $\hat{\pi}_h$ be the boosted policy obtained from Layer Boosting (Algo-
 743 rithm 7).

$$\sup_{\tilde{\pi}_h \in \Pi} \mathbb{E}_{\rho_h} \left[\sum_a Q(x, a; \pi_{h+1:H})(\tilde{\pi}_h(a | x) - \hat{\pi}_h(a | x)) \right] \leq \frac{2|\mathcal{A}|R_{\max}}{\alpha} (\varepsilon_w + \frac{2}{\sqrt{N}})$$

Algorithm 7 Layer Boosting

- 1: **input:** Sample state distribution ρ_h , policy $\pi_{h+1:H}$, number of iterations N , number of episodes M .
- 2: Initialize $\tilde{\pi}_0$ arbitrarily.
- 3: **for** $n = 1, \dots, N$ **do**
- 4: Execute $\pi_{h+1:H}$ with the sample distribution ρ_h via Algorithm 8 for M episodes, and obtain the dataset $D_n = \{(s_i, \widehat{Q}_i)\}_{i=1}^M$.
- 5: Modify dataset D_n to a new dataset $D'_n = \{(s_i, f_i)\}_{i=1}^M$, such that for all $i \in [M]$:
- 6:

$$f_i = \frac{1}{\beta} \left(y_i - \tilde{\pi}_n(\cdot | s_i) \right),$$

$$y_i = \operatorname{argmin}_{y \in \mathbb{R}^{|\mathcal{A}|}} \left\{ -\widehat{Q}_i^\top y + G \min_{z \in \Delta_{\mathcal{A}}} \|z - y\| + \frac{1}{2\beta} \left\| \tilde{\pi}_n(\cdot | s_i) - y \right\|^2 \right\},$$

where $G = H|\mathcal{A}|$, $\beta = \sqrt{\frac{1}{\alpha N}}$, and $f_i, \widehat{Q}_i \in \mathbb{R}^{|\mathcal{A}|}$.

- 7: Let \mathcal{W}_n be the policy chosen by the weak learning oracle when given data set D_n .
 - 8: Update $\tilde{\pi}_n = (1 - \eta_{2,n}) \tilde{\pi}_{n-1} + \frac{\eta_{2,n}}{\alpha} \mathcal{W}_n$.
 - 9: **end for**
 - 10: **return** $\Gamma(\tilde{\pi}_N)$.
-

Algorithm 8 Trajectory sampler

- 1: **input:** Sample state distribution ρ_h , policy $\pi_{h+1:H}$.
 - 2: Sample $s_h \sim \rho_h$, and action $a' \sim \mathcal{U}(\mathcal{A})$ uniformly.
 - 3: Take action a' at state s_h , then continue to execute $\pi_{h+1:H}$. Upon terminating, set $R(s_h, a)$ as the sum of the reward from layer $h + 1$ onwards.
 - 4: Define the vector $\widehat{Q}_{s_h}^\pi$, such that for all $a \in \mathcal{A}$, $\widehat{Q}_{s_h}^\pi(a) = |\mathcal{A}| \cdot R(s_h, a') \cdot \mathbb{I}_{a=a'}$.
 - 5: **return** $(s_h, \widehat{Q}_{s_h}^\pi)$.
-

744 The final output of PSDP Boosting is then a concatenation of the policies from layer 1 to layer H ,
 745 i.e., $\pi = (\pi_1, \pi_2, \dots, \pi_H)$. These policy from Algorithm 6 will be aggregated once again to form a
 746 final policy in Algorithm 5. Next, we provide the following result for Algorithm 5.

747 **Theorem H.2.** Algorithm 5 samples TNM episodes with $T = O\left(\frac{H^2 C_\infty^2 R_{\max}}{\varepsilon}\right)$, $N =$
 748 $O\left(\frac{HC_\infty |\mathcal{A}| R_{\max}}{\alpha \varepsilon}\right)^2$, $M = m\left(\frac{\alpha \varepsilon}{C_\infty H |\mathcal{A}| R_{\max}}, \frac{\delta}{NT}\right)$, set $\eta_{1,t} = \min\left\{1, \frac{2C_\infty}{t}\right\}$, with probability
 749 $1 - \delta$, $V^{\pi^*} - V^\pi \leq HC_\infty \mathcal{E}(\Pi, \Pi) + \varepsilon$.

750 **Sample Complexities:** If $m(\varepsilon, \delta) = \frac{\log |\mathcal{W}|}{\varepsilon^2} \log(\frac{1}{\delta})$ for some measure of weak learning complexity
751 $|\mathcal{W}|$, then the algorithm samples $O\left(\frac{H^6 C_\infty^6 R_{\max}^5 \log(|\mathcal{W}|)}{\alpha^4 \varepsilon^5}\right)$ episodes.

752 **Remark H.3.** For the discounted MDPs, we can set the effective horizon $H \approx \frac{1}{1-\gamma}$. Note that the total
753 reward R_{\max} is at most $\frac{1}{1-\gamma}$ we recover the sample complexity $\tilde{O}\left(\frac{C_\infty^6 |A|^4 \log |\mathcal{W}|}{(1-\gamma)^{11} \alpha^4 \varepsilon^5}\right)$ in the episodic
754 model (Brukhim et al., 2022, Theorem 7).

755 I Policy Cover

756 In this section, we show how a policy cover can be used to bound the occupancy mismatch term
757 C_∞ in Theorem H.2. The key idea is to construct, via the boosting procedure in Appendix H, a set
758 of policies that collectively achieve a good coverage over the state space. This guarantees that the
759 induced occupancy measure adequately covers the state space, thereby bounding the mismatch ratio
760 C_∞ .

761 A policy cover is a collection of policies whose occupancy measures together cover the state space
762 sufficiently well. Policy cover are defined over an ensemble of policies, i.e., a distribution $p \in \Delta(\bar{\Pi})$,
763 for a class of interest policy $\bar{\Pi} \subset \Pi_{\text{rns}}$. The goal is to find a distribution p such that the induced
764 occupancy measure $d_h^p(x) = \mathbb{E}_{\pi \sim p}[d_h^\pi(x)]$ covers the state space well for all $h \in [H]$. Prior analyses
765 of policy covers (Misra et al., 2020; Amortila et al., 2024) rely on such ensembles, where the expected
766 occupancy measure can be decoupled across policies, e.g., let $p = \text{Unif}(\pi^1, \dots, \pi^T)$ be the ensemble
767 then the occupancy $d_h^p(x)$ can be decoupled as $d_h^p(x) = \frac{1}{T} \sum_{t=1}^T d_h^{\pi^t}(x)$.

768 Our analysis relies on the extended MDPs and the truncated policy class, the construction is the same
769 as discounted MDP, with some modification in the context of episodic setting.

770 **Extended MDP and truncated policy class.** Similarly to the discounted MDP, we construct an
771 extended MDP $\bar{\mathcal{W}}$ by augmenting the original state space and action space with $\bar{\mathcal{A}} = \mathcal{A} \cup \{t\}$ and
772 $\bar{\mathcal{X}} = \mathcal{X} \cup \{t\}$. The transition in the extended MDP is the same as the original MDP, except that the
773 terminal state is a self-looping state, i.e., $P_h(t | t, a) = 1 \forall a \in \bar{\mathcal{A}}$. For every policy π in the original
774 class Π , we define a corresponding truncated policy $\bar{\pi} \in \bar{\Pi}$ as:

$$\bar{\pi}(x_h) = \begin{cases} \pi(x_h), & x_h \in \mathcal{B} \\ t, & \text{otherwise,} \end{cases}$$

775 for some “well-explore” set \mathcal{B} that we will defined later depend on which notion of the coverability
776 that we used. Intuitively, policy acts the same when in well-explored set, otherwise, it picks the
777 terminated action t .

778 I.1 RL Boosting with L_∞ coverability

779 We proceed with the assumption on the access to a distribution μ such that the L_∞ concentration is
780 bounded, and subsequently define our bonus exploration and the policy transformation.

781 **Assumption I.1** (L_∞ -concentrability coefficient). We assumes access to a distribution $\mu \in \Delta(\mathcal{X} \times \mathcal{A})$
782 for which the L_∞ -concentrability coefficient $\tilde{C}_\infty(\mu) := \sup_\pi \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \left\{ \frac{d_h^\pi(x,a)}{\mu(x,a)} \right\}$ is bounded.

783 We abbreviate $C_\infty \equiv C_{\infty;h}^M(\mu)$, and define our exploration bonus as:

$$\begin{aligned} w_h^k(x, a) &= \frac{\tilde{C}_\infty \mu(x, a)}{\sum_{i < t} d_h^{\pi^i}(x, a) + \tilde{C}_\infty \mu(x, a)}, \\ b_h^k(x, a) &= 2w_h^k(x, a). \end{aligned} \tag{23}$$

Policy transformation $\mathcal{T}(\Pi, b^k)$. Given the well-explored set $\mathcal{B}^k := \{x_h : w_h^k(x, a) \leq \frac{1}{2}, \forall a \in \mathcal{A}, h \in [H]\}$, the transformed policy $\hat{\pi}$ is defined as follows:

$$\hat{\pi}(x) = \begin{cases} \pi(x_h), & \text{if } x_h \in \mathcal{B}^k \\ \text{Unif}\left(\{a \in \mathcal{A} : w_h^k(x_h, a) > \frac{1}{2}\}\right), & \text{otherwise} \end{cases}$$

Algorithm 9 Boosting with L_∞ coverability

- 1: **input:** number of epoch K , number of iterations T , number of weak learner N , number of sample episodes M , policy class Π .
 - 2: Initialize policy π^1 arbitrarily.
 - 3: **for** $k = 1, \dots, K$ **do**
 - 4: Construct new bonus $b_h^k(x, a)$ with policy cover $d_h^{p^k}$ (Eq. (23)).
 - 5: Update policy class $\widehat{\Pi} \leftarrow \mathcal{T}(\Pi, b^k)$.
 - 6: Run (Algorithm 6) with reward function $r + b^k$, policy class $\widehat{\Pi}$, initial distribution d^{p^k} to obtain π^k .
 - 7: Update the policy cover $p^{k+1} = \text{Unif}(\pi_1, \dots, \pi_k)$.
 - 8: **end for**
 - 9: **return** $\pi := \text{argmax}_{\pi \in \{\pi^1, \dots, \pi^K\}} V^\pi$.
-

784 With the definition of the bonus function and the policy transformation, we provide the boosting for
 785 episodic MDP with the L_∞ coverability as follows.

Theorem I.2. For any $\varepsilon, \delta \in (0, 1)$, suppose Assumption I.1 is satisfied, Algorithm 9 is instantiated with the following parameters: $K = O(\frac{\widetilde{C}_\infty H \iota}{\varepsilon})$, $N = O(\frac{\widetilde{C}_\infty^2 H^4 |\mathcal{A}|^2 \iota^2}{\alpha^2 \varepsilon^4})$, $M = m\left(\frac{\alpha \varepsilon^2}{\widetilde{C}_\infty H^2 |\mathcal{A}|}, \frac{\delta}{KHN}\right)$, where $\iota = \log(\frac{\widetilde{C}_\infty H}{\varepsilon})$. Then Algorithm 2 produces policy π , such that with probability at least $1 - \delta$,

$$V^* - V^\pi \leq \frac{\widetilde{C}_\infty H \iota}{\varepsilon} \mathcal{E}(\Pi, \Pi) + \varepsilon.$$

786 **Sample Complexity:** If $m(\varepsilon, \delta) = \frac{\log |\mathcal{W}|}{\varepsilon^2} \log \frac{1}{\delta}$ for some measure of weak learning complexity $|\mathcal{W}|$,
 787 then the algorithm use at most $KHN M \leq \widetilde{O}\left(\frac{\widetilde{C}_\infty^5 H^{10} |\mathcal{A}|^4 \log |\mathcal{W}|}{\alpha^4 \varepsilon^9}\right)$ episodes.

788 *Proof.* For any episode $k \in [K]$, we define a well-explored set $\mathcal{B}_h^k := \{x, a : \frac{\widetilde{C}_\infty \mu(x, a)}{d^{p^k}(x, a)} \leq k - 1\}$,
 789 we say $x \in \mathcal{B}_h^k$ if $(x, a) \in \mathcal{B}_h^k, \forall a \in \mathcal{A}$. For any policy $\bar{\pi} \in \bar{\Pi}_k$, and $\pi \in \mathcal{T}(\Pi, \mathcal{B}^k)$, for any $x \in \mathcal{B}^k$,
 790 we have

$$d^{\bar{\pi}}(x) \leq (k - 1) d^{p^k}(x)$$

$$d_h^{\bar{\pi}}(x_h) = \sum_{a \in \mathcal{A}} d_h^{\bar{\pi}}(x_h, a) \leq \sum_{a \in \mathcal{A}} \widetilde{C}_\infty \mu(x_h, a) \leq \sum_{a \in \mathcal{A}} (k - 1) d^{p^k}(x_h, a) \leq (k - 1) d^{p^k}(x_h),$$

791 Let $\pi = \text{argmin}_{\pi \in \{\pi^1, \dots, \pi^K\}} V^{\pi^k}$ be the output of the Algorithm 9, we have the following:

$$\leq \frac{1}{K} \sum_{k=1}^K V^{\pi^*} - V^{\pi^k}$$

792 We proceed to bound the RHS by decomposing the summation into the following three terms:

$$\frac{1}{K} \sum_{k=1}^K V^{\pi^*} - V^{\pi^k} = \underbrace{\frac{1}{K} \sum_{k=1}^K V^{\pi^*} - \overline{V}_{b^k}^{\pi^*}}_{I_1 \leq 0} + \underbrace{\frac{1}{K} \sum_{k=1}^K \overline{V}_{b^k}^{\pi^*} - \overline{V}_{b^k}^{\pi^k}}_{I_2} + \underbrace{\frac{1}{K} \sum_{k=1}^K \overline{V}_{b^k}^{\pi^k} - V^{\pi^k}}_{I_3}.$$

793 For I_2 , we have

$$\begin{aligned} \bar{V}_{b^k}^{\bar{\pi}^k} - \bar{V}_{b^k}^{\pi^k} &= \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\bar{Q}_h \left(x_h; \bar{\pi}_h \circ \pi_{h+1:H}^k \right) - \bar{V}_h \left(x_h; \pi_{h:H}^k \right) \right] \\ &= \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\bar{Q}_h \left(x_h; \bar{\pi}_h \circ \pi_{h+1:H}^k \right) - \bar{V}_h \left(x_h; \pi_{h:H}^k \right) \mathbb{1}\{x \in \mathcal{B}^k\} \right] \\ &\quad + \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\bar{Q}_h \left(x_h; \bar{\pi}_h \circ \pi_{h+1:H}^k \right) - \bar{V}_h \left(x_h; \pi_{h:H}^k \right) \mathbb{1}\{x \notin \mathcal{B}^k\} \right] \end{aligned}$$

794 The first term in the equation above can be upper bounded by:

$$\begin{aligned} &\sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\bar{Q}_h \left(x_h; \bar{\pi}_h \circ \pi_{h+1:H}^k \right) - \bar{V}_h \left(x_h; \pi_{h:H}^k \right) \right] \\ &\leq (k-1) \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{b^k}} \left[\bar{Q}_h \left(x_h; \bar{\pi}_h \circ \pi_{h+1:H}^k \right) - \bar{V}_h \left(x_h; \pi_{h:H}^k \right) \right] \\ &\leq (k-1) \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{b^k}} \left[\bar{Q}_h \left(x_h; \bar{\pi}_h \circ \pi_{h+1:H}^k \right) - \max_a \bar{Q}_h(x_h, a; \bar{\pi}_h \circ \pi_{h+1:H}^k) \right] \\ &\quad + (k-1) \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{b^k}} \left[\max_a \bar{Q}_h(x_h, a; \pi_{h+1:H}^k) - \max_{\pi_h \in \Pi} \bar{Q}_h(x_h, a; \pi_h \circ \pi_{h+1:H}^k) \right] \\ &\quad + (k-1) \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{b^k}} \left[\max_{\pi_h \in \Pi} \bar{Q}_h(x_h, a; \pi_h \circ \pi_{h+1:H}^k) - \bar{V}_h \left(x_h; \pi_{h:H}^k \right) \right] \\ &\leq (k-1) H \mathcal{E}(\Pi, \Pi) + \frac{2(k-1)H|\mathcal{A}|}{\alpha} (\varepsilon_w + \frac{2}{\sqrt{N}}), \end{aligned}$$

795 where the last inequality follows by Lemma H.1.

796 For the second term, we have $\bar{\pi}_h$ plays t when $x_h \notin \mathcal{B}^k$, this implies $\bar{Q}_h \left(x_h; \bar{\pi}_h \circ \pi_{h+1:H}^k \right) = 1$.

797 By the policy transformation construction, we have $\bar{V}_h \left(x_h; \pi_{h:H}^k \right) \geq 1$, therefore, we can safely
798 ignore this second term.

799 For I_3 , we have:

$$\begin{aligned} \sum_{k=1}^K \bar{V}_{b^k}^{\pi^k} - V^{\pi^k} &= \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}^{\pi^k} [b_h^k(x, a)] \\ &= \sum_{h=1}^H \sum_{k=1}^K \mathbb{E}^{\pi^k} \left[\frac{\tilde{C}_\infty \mu(x_h, a_h)}{\sum_{i < k} d_h^{\pi^i}(x) + \tilde{C}_\infty \mu(x_h, a_h)} \right] \\ &= 2\tilde{C}_\infty \sum_{x \in \mathcal{X}, a \in \mathcal{A}} \sum_{k=1}^K \mu(x, a) \frac{d_h^{\pi^k}(x, a)}{\sum_{i < k} d_h^{\pi^i}(x, a) + \tilde{C}_\infty \mu(x, a)} \\ &\leq 2H\tilde{C}_\infty \log(2K) \end{aligned}$$

800 Put everything together, we have:

$$V^{\pi^*} - V^{\pi} \leq KH \mathcal{E}(\Pi, \Pi) + \frac{2KH|\mathcal{A}|}{\alpha} (\varepsilon_w + \frac{2}{\sqrt{N}}) + \frac{2H\tilde{C}_\infty \log(2K)}{K}$$

801

□

802 I.2 RL boosting with pushforward coverability

803 The exploration bonus that uses L_∞ does require the knowledge of the reference distribution μ . The
804 pushforward coverability aim to remove this assumption, however, we would get a looser upper bound.

805 This pushforward coverability is a more general notion that can be applied to various structured
 806 MDPs. The main idea is to learn a new policy that can reach the states that are not well covered by
 807 the previous policy cover. Instead of using a reference distribution μ in L_∞ coverability, we use the
 808 transition instead. Specifically, we use the following weight function as an exploration objective:

$$w_h^k(x_h | x_{h-1}, a_{h-1}) := \frac{P_{h-1}(x_h | x_{h-1}, a_{h-1})}{\sum_{i < k} d_h^{\pi^{h,i}}(x_h) + P_{h-1}(x_h | x_{h-1}, a_{h-1})},$$

809 Intuitively, if a state x_h is poorly covered by previous policies, the sum $\sum_{i < k} d_h^{\pi^{h,i}}(x_h)$ will be small,
 810 and the resulting weight will be higher, thereby encouraging exploration of that state. Similar to L_∞
 811 coverability, we use the following exploration bonus: $b_h^k(x_h | x_{h-1}, a_{h-1}) := w_h^k(x_h | x_{h-1}, a_{h-1})$

812 Since the reward construction requires the knowledge of the transition, which is unknown in general,
 813 we can efficiently estimate the weight function via contrastive learning (Amortila et al., 2024,
 814 Algorithm 4).

815 Next, we start with the definition of pushforward coverability as follow

Definition I.3 (Pushforward coverability).

$$C_{\text{push};h} = \inf_{\mu \in \Delta(\mathcal{X})} \sup_{(x,a,x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}} \left\{ \frac{P_{h-1}(x' | x, a)}{\mu(x')} \right\},$$

816 with $C_{\text{push}} = \max_{h \in [H]} C_{\text{push};h}$.

817 In contrast to the distribution mismatch C_∞ , the pushforward coverability C_{push} is bounded by the
 818 intrinsic complexity of the MDP; this is general notion that can be applied to various structured
 819 MDPs. For example, tabular MDPs with state space \mathcal{X} have $C_{\text{push}} \leq |\mathcal{X}|$, block MDPs with latent
 820 space \mathcal{S} have $C_{\text{push}} \leq |\mathcal{S}|$, and low-rank MDPs of dimension d have $C_{\text{push}} \leq d$.

Algorithm 10 Boosting with push forward coverability

- 1: **input:** number of iterations T , number of weak learner N , number of sample episodes M , policy class Π .
 - 2: Initialize policy π^1 arbitrarily, $p^1 = \{\pi^1\}$.
 - 3: **for** $k = 1, \dots, K$ **do**
 - 4: Run Algorithm 11 to estimate $\hat{w}_h^k(x | x', a')$ with policy cover p^k and then construct new bonus $\hat{b}_h^k(x | x', a')$ for each layer h .
 - 5: Run Algorithm 6 with reward function $r + \hat{b}_h^k$, initial distribution d^{p^k} to obtain π^k .
 - 6: Update the policy cover $p^{k+1} = \text{Unif}(\pi_1, \dots, \pi_k)$.
 - 7: **end for**
 - 8: **return** $\pi := \arg\max_{\pi \in \{\pi^1, \dots, \pi^K\}} V^\pi$.
-

821 To perform the weight estimation, we assume access to a weight function class $\mathcal{W} = \mathcal{W}_{1:H}$, with
 822 $\mathcal{W}_h \subseteq (\mathcal{X} \times \mathcal{A} \times \mathcal{X} \rightarrow \mathbb{R}_+)$ which we restate the following assumption from Amortila et al. (2024)
 823 as follows.

824 **Assumption I.4** (Weight function realizability). For all $h \geq 2$ and all $\pi \in \Pi$

$$w_h^\pi(x' | x, a) := \frac{P_{h-1}(x' | x, a)}{d_h^\pi(x')} \in \overline{\mathcal{W}}_h$$

825

826 Next, we provide the result for Algorithm 10 in the following theorem.

Theorem I.5. For any $\varepsilon, \delta \in (0, 1)$, suppose Assumption I.4 is satisfied. Algorithm 10 is instantiated with the following parameters: $K = O\left(\frac{C_{\text{push}} H^\iota}{\varepsilon}\right)$, $N = O\left(\frac{C_{\text{push}}^2 H^6 |\mathcal{A}|^2 \iota^2}{\alpha^2 \varepsilon^6}\right)$, $M = m\left(\frac{\alpha \varepsilon^3}{C_{\text{push}} H^3 |\mathcal{A}|}, \frac{\delta}{K H N N_{\text{weight}}}\right)$, $N_{\text{weight}} = O\left(\frac{C_{\text{push}} H^4 |\mathcal{A}| \log(|\overline{\mathcal{W}}| \delta^{-1}) \iota}{\varepsilon^4}\right)$, where $\iota = \log\left(\frac{C_{\text{push}} H}{\varepsilon}\right)$. Then Algorithm 10 produces policy π , such that with probability at least $1 - \delta$,

$$V^* - V^\pi \leq \frac{C_{\text{push}} H^3 \iota}{\varepsilon^2} \mathcal{E}(\Pi, \Pi) + \varepsilon.$$

Algorithm 11 Weight estimation

- 1: Let $\frac{1}{(t-1)} \sum_{i < t} \pi^i \circ_{h-1} \pi_{\text{unif}}$ if $t \geq 1$ and $q := p_{h-1}$ otherwise.
 - 2: Let $\mathcal{D}_1 = \mathcal{D}_2 = \emptyset$.
 - 3: For each $j \in [n]$, draw $\pi \sim q$ and sample $(x_{h-1}^j, a_{h-1}^j, x_h^j) \sim \pi$. Add $(x_{h-1}^j, a_{h-1}^j, x_h^j)$ to both \mathcal{D}_1 and \mathcal{D}_2 .
 - 4: **for** $i < t$ **do**
 - 5: Draw n samples $\left\{ (x_{h-1}^j, a_{h-1}^j, x_h^j) \right\}_{j \in [n]}$ independently by drawing $\pi \sim q$ and $(x_{h-1}^j, a_{h-1}^j, x_h^j) \sim \pi$.
 - 6: Draw n samples $\left\{ (\tilde{x}_{h-1}^j, \tilde{a}_{h-1}^j, \tilde{x}_h^j) \right\}_{j \in [n]}$ by sampling $(\tilde{x}_{h-1}^j, \tilde{a}_{h-1}^j, \tilde{x}_h^j) \sim \pi^i$.
 - 7: Add $\left\{ (x_{h-1}^j, a_{h-1}^j, x_h^j) \right\}_{j \in [n]}$ to \mathcal{D}_1 and add $\left\{ (x_{h-1}^j, a_{h-1}^j, \tilde{x}_h^j) \right\}_{j \in [n]}$ to \mathcal{D}_2 .
 - 8: **end for**
 - 9: **Set** $\hat{w} := \arg \max_{w \in \mathcal{W}_h} \hat{\mathbb{E}}_{\mathcal{D}_1} \left[\log \left(w(x_h | x_{h-1}, a_{h-1}) \right) \right] - t \cdot \hat{\mathbb{E}}_{\mathcal{D}_2} \left[w(x_h | x_{h-1}, a_{h-1}) \right]$.
-

827 Assuming access to an efficient weak learning oracle, our result implies that both computational and
 828 statistical complexities scale as $\text{poly}(C_{\text{push}}, H, |\mathcal{A}|, \log |\mathcal{W}|, \log(|\bar{\mathcal{W}}|) \alpha^{-1}, \varepsilon^{-1})$.

829 Specifically,

830 *Proof of Theorem I.5.* We follow the same proof procedure in L_∞ cover, we have the following:

$$V^{\pi^*} - V^\pi \leq \frac{1}{K} \sum_{k=1}^K V^{\pi^*} - V^{\pi^k}$$

831 We proceed to bound the RHS by decomposing the summation into the following three terms:

$$\sum_{k=1}^K V^{\pi^*} - V^{\pi^k} = \underbrace{\frac{1}{K} \sum_{k=1}^K V^{\pi^*} - \bar{V}_{\hat{b}^k}}_{I_1} + \underbrace{\frac{1}{K} \sum_{k=1}^K \bar{V}_{\hat{b}^k} - \bar{V}_{\hat{b}^k}}_{I_2} + \underbrace{\frac{1}{K} \sum_{k=1}^K \bar{V}_{\hat{b}^k} - V^{\pi^k}}_{I_3},$$

832 We have $I_1 \leq 0$ by the dynamics of the extended MDP.

833 For I_2 with $x \in \mathcal{B}^k$ can be upper bounded by $\frac{k-1}{\zeta} H \mathcal{E}(\Pi, \Pi) + \frac{(k-1)H|\mathcal{A}|}{\zeta \alpha} (\varepsilon_w + \frac{2}{\sqrt{N}})$, follow by a
 834 similar argument in L_∞ .

835 For I_2 with $x \notin \mathcal{B}^k$

$$\begin{aligned} & \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\bar{Q}_{\hat{b}^k} \left(x_h; \bar{\pi}_h \circ \pi_{h+1:H}^k \right) - \bar{V}_{\hat{b}^k} \left(x_h; \pi_{h:H}^k \right) \mathbb{1}\{x \notin \mathcal{B}^k, x \neq \mathbf{t}\} \right] \\ & \leq \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[1 - 2\mathbb{E}_{x', a'}^{\pi} [\hat{w}_h^k(x | x', a')] \mathbb{1}\{x \notin \mathcal{B}^k, x \neq \mathbf{t}\} \right] \\ & \leq \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[1 - 2\mathbb{E}_{x', a'}^{\pi} [w_h^k(x | x', a')] \mathbb{1}\{x \notin \mathcal{B}^k, x \neq \mathbf{t}\} \right] \\ & \quad + 2\mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\mathbb{E}_{x', a'}^{\pi} [w_h^k(x | x', a')] - \hat{w}_h^k(x | x', a')] \mathbb{1}\{x \notin \mathcal{B}^k, x \neq \mathbf{t}\} \right] \end{aligned}$$

836 We further bound the first term as follows:

$$\sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\sum_{x', a' \in \mathcal{X} \times \mathcal{A}} d^\pi(x', a') P_{h-1}(x_h | x', a') \hat{w}_h^k(x_h | x', a') \mathbb{1}\{x \notin \mathcal{B}^k, x \neq \mathbf{t}, P_{h-1}(x_h | x', a') < \zeta\} \right] < H\zeta$$

837 For the second term,

$$\begin{aligned}
& \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\mathbb{E}_{x_{h-1}, a_{h-1}}^{\pi} [w_h^k(x_h | x_{h-1}, a_{h-1}) - \widehat{w}_h^k(x_h | x_{h-1}, a_{h-1})] \mathbb{1}\{x_h \notin \mathcal{B}^k, x_h \neq \mathfrak{t}\} \right] \\
& \leq 2 \mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \left[\sqrt{\mathbb{E}_{x_{h-1}, a_{h-1}}^{\pi} \left[\left(\sqrt{w_h^k(x_h | x_{h-1}, a_{h-1})} - \sqrt{\widehat{w}_h^k(x_h | x_{h-1}, a_{h-1})} \right)^2 \mathbb{1}\{x_h \notin \mathcal{B}^k, x_h \neq \mathfrak{t}\} \right]} \right] \\
& \leq 2 \sqrt{\mathbb{E}_{x_h \sim d_h^{\bar{\pi}}} \mathbb{E}_{x_{h-1}, a_{h-1}}^{\pi} \left[\left(\sqrt{w_h^k(x_h | x_{h-1}, a_{h-1})} - \sqrt{\widehat{w}_h^k(x_h | x_{h-1}, a_{h-1})} \right)^2 \mathbb{1}\{x_h \notin \mathcal{B}^k, x_h \neq \mathfrak{t}\} \right]} \\
& \leq 2 \sqrt{\sum_{\substack{x_h \in \mathcal{X}: \\ x_h \neq \mathfrak{t}}} d_h^{\bar{\pi}}(x_h) \sum_{\substack{x', a' \in \mathcal{X} \times \mathcal{A}: \\ x' \neq \mathfrak{t}, a' \neq \mathfrak{t}}} d^{\pi}(x') \pi(a' | x') P_{h-1}(x_h | x', a') \left(\sqrt{w_h^k(x_h | x', a')} - \sqrt{\widehat{w}_h^k(x_h | x', a')} \right)^2} \\
& \leq 2 \sqrt{\sum_{\substack{x_h \in \mathcal{X}: \\ x_h \neq \mathfrak{t}}} d_h^{\bar{\pi}}(x_h) \sum_{\substack{x', a' \in \mathcal{X} \times \mathcal{A}: \\ x' \neq \mathfrak{t}, a' \neq \mathfrak{t}}} (k-1) d^{p^k}(x') \pi(a' | x') P_{h-1}(x_h | x', a') \left(\sqrt{w_h^k(x_h | x', a')} - \sqrt{\widehat{w}_h^k(x_h | x', a')} \right)^2} \\
& \leq 2 \sqrt{|\mathcal{A}| \sum_{\substack{x_h \in \mathcal{X}: \\ x_h \neq \mathfrak{t}}} \sum_{\substack{x', a' \in \mathcal{X} \times \mathcal{A}: \\ x' \neq \mathfrak{t}, a' \neq \mathfrak{t}}} (k-1) d^{p^k}(x') \pi_{\text{unif}}(a' | x') P_{h-1}(x_h | x', a') \left(\sqrt{w_h^k(x_h | x', a')} - \sqrt{\widehat{w}_h^k(x_h | x', a')} \right)^2} \\
& \leq 2 \sqrt{\frac{(k-1)|\mathcal{A}|}{\zeta(k-1)} \sum_{i < t} \mathbb{E}^{\pi^i \circ_{h-1} \pi_{\text{unif}}} \left[\left(\sqrt{w_h^k(x_h | x', a')} - \sqrt{\widehat{w}_h^k(x_h | x', a')} \right)^2 \right]} \\
& \leq \sqrt{\frac{(k-1)|\mathcal{A}| \varepsilon_{\text{weight}}^2}{\zeta}}
\end{aligned}$$

838 We have

$$I_2 \leq \frac{k-1}{\zeta} H \mathcal{E}(\Pi, \Pi) + \frac{(k-1)H|\mathcal{A}|}{\zeta \alpha} (\varepsilon_w + \frac{2}{\sqrt{N}}) + \sqrt{\frac{(k-1)|\mathcal{A}| \varepsilon_{\text{weight}}^2}{\zeta}} + H \zeta$$

839 For I_3 ,

$$\begin{aligned}
I_3 &= \frac{1}{K} \sum_{k=1}^K \mathbb{E}^{\pi^k} [\widehat{w}_h^k(x_h | x_{h-1}, x_{h-1})] \\
&\leq \frac{1}{K} \sum_{h=1}^H \sum_{k=1}^K \left(3 \mathbb{E}_h^{\pi^k} [w_h^k(x_h | x_{h-1}, x_{h-1})] + 2 \mathbb{E}_h^{\pi^k} \left[\left(\sqrt{\widehat{w}_h^k(x_h | x_{h-1}, x_{h-1})} - \sqrt{w_h^k(x_h | x_{h-1}, x_{h-1})} \right)^2 \right] \right)
\end{aligned}$$

840 Applying Lemma L.6, we have

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E}_h^{\pi^k} [(\sqrt{\widehat{w}_h^k} - \sqrt{w_h^k})^2] \\
& \leq \sqrt{2C_{\text{push};h} \log(K) \sum_{\substack{k \in [K] \\ i < k}} \mathbb{E}^{\pi^i} \left[(\mathbb{E}[(\sqrt{\widehat{w}_h^k}(x_h | x_{h-1}, a_{h-1}) - \sqrt{w_h^k}(x_h | x_{h-1}, a_{h-1}))^2 | x_{h-1}, a_{h-1} = \pi^k(x_{h-1})]) \right]} \\
& \quad + 4C_{\text{push};h} \\
& \leq \sqrt{8C_{\text{push};h} |\mathcal{A}| \log(K) \sum_{k=1}^K \sum_{i < k} \mathbb{E}^{\pi^i \circ_{h-1} \pi_{\text{unif}}} \left[\left(\sqrt{\widehat{w}_h^k}(x_h | x_{h-1}, a_{h-1}) - \sqrt{w_h^k}(x_h | x_{h-1}, a_{h-1}) \right)^2 \right]} + 4C_{\text{push};h} \\
& \leq \sqrt{8|\mathcal{A}| C_{\text{push};h} \log(K) \varepsilon_{\text{weight}}^2} + 4C_{\text{push}},
\end{aligned}$$

841 where we use Lemma E.5 in the last inequality. The bound relate to the weight can be control by the
 842 following lemma:

Lemma I.6.

$$\sum_{k=1}^K \mathbb{E}_h^{\pi^k} \left[w_h^k(x_h | x_{h-1}, x_{h-1}) \right] \leq C_{\text{push}} \log(2K).$$

843 Therefore, we can bound I_3 as follows,

$$I_3 \leq \frac{4HC_{\text{push}} \log(2K)}{K} + H \sqrt{8|\mathcal{A}|C_{\text{push}} \log(K)\varepsilon_{\text{weight}}^2} + \frac{4HC_{\text{push}}}{K},$$

844 Put everything together, we have

$$\begin{aligned} V^{\pi^*} - V^{\pi} &\leq \frac{KH}{\zeta} \mathcal{E}(\mathbb{P}, \Pi) + \frac{KH|\mathcal{A}|}{\zeta\alpha} (\varepsilon_w + \frac{2}{\sqrt{N}}) + H \sqrt{\zeta^{-1}K|\mathcal{A}|\varepsilon_{\text{weight}}^2} \\ &\quad + \frac{8HC_{\text{push}} \log(2K)}{K} + H \sqrt{8|\mathcal{A}|C_{\text{push}} \log(K)\varepsilon_{\text{weight}}^2} + H\zeta \end{aligned}$$

845 We complete the proof by choosing $\zeta = H^{-1}\varepsilon$. □

846 *Proof of Lemma I.6.* Let $\widetilde{d}_h^k(x_h) = \sum_{i=1}^{k-1} d_h^{\pi^i}(x_h)$, for any $h \in [H]$, let $\mu \in \Delta(\mathcal{X})$ attain the value
 847 of $C_{\text{push};h}$, with $\varepsilon = 1$, and $\delta = C_{\text{push};h}$, applying Lemma L.4

$$\begin{aligned} \mathbb{E}^{\pi^k} \left[\frac{P_{h-1}(x_h | x_{h-1}, a_{h-1})}{\widetilde{d}_h^k(x_h) + P(x_h | x_{h-1}, a_{h-1})} \right] &\leq \mathbb{E}^{\pi^k} \left[\frac{P_{h-1}(x_h | x_{h-1}, a_{h-1})}{\widetilde{d}_h^k(x_h) + C_{\text{push};h}^M \mu(x_h)} \right] \\ &\quad + C_{\text{push};h} \cdot \mathbb{E}^{\pi^k} \left[\frac{\mu(x_h)}{\widetilde{d}_h^k(x_h) + C_{\text{push};h} \mu(x_h)} \right] \\ &\leq 2C_{\text{push};h} \cdot \mathbb{E}^{\pi^k} \left[\frac{\mu(x_h)}{\widetilde{d}_h^k(x_h) + C_{\text{push};h}^M \mu(x_h)} \right]. \end{aligned}$$

848 Since $\sup_{\pi \in \Pi} d_h^{\pi}(x) \leq \sup_{x' \in \mathcal{X}, a \in \mathcal{A}} P_{h-1}(x | x', a) \leq C_{\text{push};h}^M \mu(x)$ for all $x \in \mathcal{X}$, Lemma L.3
 849 implies that

$$\sum_{x \in \mathcal{X}} \sum_{k=1}^K \mu(x) \frac{d^{\pi^k}(x)}{\widetilde{d}_h^k(x) + C_{\text{push};h} \mu(x)} \leq 2 \log(2K),$$

850 which concludes the proof. □

851 I.3 Weight Function Estimation.

852 **Lemma I.7.** For any $k \in [K]$, for any $\varepsilon_{\text{weight}}, \delta \in (0, 1)$, distribution $p^k \in \Delta(\mathbb{P})$ and $\pi^1, \dots, \pi^1 \in$
 853 \mathbb{P} , Algorithm 11 ensures that with probability at least $1 - \delta$, the output \widehat{w}_h^k satisfies

$$\frac{1}{k-1} \sum_{i < t} \mathbb{E}^{M^*, \pi_h^i \circ_{h-1} \pi_{\text{unif}}} \left[\left(\sqrt{\widehat{w}_h^i(x_h | x_{h-1}, a_{h-1})} - \sqrt{w_h^i(x_h | x_{h-1}, a_{h-1})} \right)^2 \right] \leq \varepsilon^2,$$

854 and does so using at most $N_{\text{weight}} = 40 \frac{\log(|\mathcal{W}|\delta^{-1})}{\varepsilon_{\text{weight}}^2}$ episodes.

855 *Proof.* This proof follows similarly to Lemma J.7 of Amortila et al. (2024), Let $\bar{w}_h^t := t \cdot w_h^t$ and let
 856 $\tilde{w}_h^t := t \cdot \widehat{w}_h^t$. Here, we solve the optimization problem over the class $t \cdot \mathcal{W}_h$, which has $\|w'\|_{\infty} \leq t$
 857 for all $w' \in \mathcal{W}_h$. Similar to the proof in Amortila et al. (2024), we have

$$\mu(x' | x, a) = P_{h-1}(x' | x, a), \quad \nu(x' | x, a) = \frac{1}{t} \left(\sum_{i < t} d_h^{M^*, \pi^{h,i}}(x') + P_{h-1}(x' | x, a) \right),$$

858 and

$$\omega(x, a) = \frac{1}{t-1} \sum_{i < t} d_{h-1} \pi^{h,i} \circ_{h-1} \pi_{\text{unif}}(x, a).$$

859 Under the weight assumption, we have

$$\frac{\mu(x' | x, a)}{\nu(x' | x, a)} = \bar{w}_h^t(x' | x, a) \in t \cdot \mathcal{W}_h,$$

860 so Lemma E.5 imply that

$$\mathbb{E}_{(x_{h-1}, a_{h-1}) \sim \omega} \left[\left(\sqrt{\bar{w}_h^t(x_h | x_{h-1}, a_{h-1})} - \sqrt{\hat{w}_h^t(x_h | x_{h-1}, a_{h-1})} \right)^2 \right] \leq \frac{20t \log(|\mathcal{W}|\delta^{-1})}{n},$$

861 or equivalently,

$$\mathbb{E}_{(x_{h-1}, a_{h-1}) \sim \omega} \left[\left(\sqrt{\hat{w}_h^t(x_h | x_{h-1}, a_{h-1})} - \sqrt{w_h^t(x_h | x_{h-1}, a_{h-1})} \right)^2 \right] \leq \frac{20 \log(|\mathcal{W}|\delta^{-1})}{n}.$$

862

□

863 I.4 Proof of Appendix H

864 *Proof of Theorem H.2.* We define $\|\cdot\|_{\infty,1} : \mathbb{R}^{|\mathcal{X}| \times |\mathcal{A}|} \rightarrow \mathbb{R}$ as $\|z\|_{\infty,1} = \max_{x \in \mathcal{X}} \sum_a |z_{x,a}|$, and
865 the following lemma give the result of the smoothness of V function for this norm

866 **Lemma I.8** (Smoothness). *For any two policies $\pi^{(1)}, \pi^{(2)}$*

$$\left| V^{\pi^{(1)}} - V^{\pi^{(2)}} - \langle \nabla_{\pi} V^{\pi^{(2)}}, \pi^{(1)} - \pi^{(2)} \rangle \right| \leq H^2 \cdot R_{\max} \left\| \pi^{(1)} - \pi^{(2)} \right\|_{\infty,1}^2,$$

867 *in other words, V^{π} is $H^2 R_{\max}$ in $\|\cdot\|_{\infty,1}$ norm*

868 Similar to [Brukhim et al. \(2022\)](#), the rest of the proof essentially following the Frank-Wolfe style
869 optimization, we apply the gradient dominance for the V function, we state such result in the next
870 lemma.

871 **Lemma I.9** (Gradient domination). $V_1(\pi^*) - V_1(\pi) \leq C_{\infty}(H\mathcal{E}(\mathbb{I}, \Pi) + \max_{\pi' \in \Pi^H} \langle \nabla_{\pi} V^{\pi}, \pi' - \pi \rangle)$

872 Remind that, the FW algorithm in [Brukhim et al. \(2022\)](#) define the boosting Oracle \mathcal{O} as an $(\epsilon_0, \mathcal{K})$ -
873 approximate linear optimizer over a convex set $\mathcal{K} \subseteq \mathbb{R}^d$ such that for any given $v \in \mathbb{R}^d$, we have
874 $v^{\top} \mathcal{O}(v) \geq \max_{u \in \mathcal{K}} v^{\top} u - \epsilon_0$. In the context of boosting RL, $v = \nabla_{\pi} V^{\pi}$, and let $\hat{\pi} \leftarrow \mathcal{O}(\nabla_{\pi} V^{\pi})$
875 be policy obtained from the Oracle. This Oracle will guarantee that $\langle \nabla_{\pi} V^{\pi}, \hat{\pi} \rangle \geq \sup_{\bar{\pi}} \langle \nabla_{\pi} V^{\pi}, \bar{\pi} \rangle -$
876 $\epsilon_{\mathcal{O}}$. Next, let $h_t = V_1(\pi^*) - V_1(\pi_t)$, we have

$$\begin{aligned} h_t &\leq h_{t-1} - \eta_t \langle \nabla_{\pi} V^{\pi_{t-1}}, \hat{\pi}_t - \pi_{t-1} \rangle + \eta_t^2 H^2 R_{\max} \|\hat{\pi}_t - \pi_{t-1}\|_{\infty,1}^2 && \text{(Lemma L.9)} \\ &\leq h_{t-1} - \sup_{\bar{\pi}} \eta_t \langle \nabla_{\pi} V^{\pi_{t-1}}, \bar{\pi} - \pi_{t-1} \rangle + \eta_t^2 H^2 R_{\max} \|\hat{\pi}_t - \pi_{t-1}\|_{\infty,1}^2 + \eta_t \epsilon_{\mathcal{O}} && \text{(Oracle guarantee)} \\ &\leq \left(1 - \frac{\eta_t}{C_{\infty}}\right) h_{t-1} + \eta_t^2 H^2 R_{\max} \|\hat{\pi}_t - \pi_{t-1}\|_{\infty,1}^2 + \eta_t (\epsilon_{\mathcal{O}} + H\mathcal{E}) && \text{(Gradient dominance, Lemma I.9)} \end{aligned}$$

(24)

877 The Oracle for our method is PSDP Boosting, let $\hat{\pi}$ be the policy obtain from this algorithm, and
 878 $\bar{\pi} = \sup_{\pi'} \langle \nabla_w V^\pi, \pi' \rangle$, we proceed to bound $\langle \nabla_\pi V^\pi, \bar{\pi} - \pi \rangle$ by applying Lemma I.10:

$$\begin{aligned} \langle \nabla_\pi V^{\pi_{t-1}}, \bar{\pi} - \pi_{t-1} \rangle &= \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\pi_{t-1}}} \left[\sum_{a \in \mathcal{A}} Q_h^\pi(x_h, a) \left(\bar{\pi}_h(a | x_h) - \pi_h(a | x_h) \right) \right] \\ &\leq \frac{2H|\mathcal{A}|R_{\max}}{\alpha} \left(\varepsilon_w + \frac{2}{\sqrt{N}} \right), \end{aligned}$$

879 where the last inequality comes from the Layer Boosting algorithm, which we state the following
 880 lemma.

881 Put everything together, we have

$$h_t \leq \left(1 - \frac{\eta_t}{C_\infty}\right) h_{t-1} + \eta_t H C_\infty \left(\mathcal{E} + \frac{2|\mathcal{A}|R_{\max}}{\alpha} \left(\varepsilon_w + \frac{2}{\sqrt{N}} \right) \right) + \eta_t^2 H^2 R_{\max}$$

882 Applying the following result of a bounded positive sequences

883 **Claim 3** (Claim 21, (Brukhim et al., 2022)). *Let $C \geq 1$, let g_t be the B -bounded positive sequence*
 884 *such that*

$$h_t \leq \left(1 - \frac{\sigma_t}{C}\right) h_{t-1} + \sigma_t E + \sigma_t^2 D,$$

885 *then choosing $\sigma_t = \min\{1, \frac{2C}{t}\}$ implies $g_t \leq \frac{2C^2 \max\{2D, B\}}{t} + E$*

886 We can bound $V^{\pi^*} - V^{\pi_T}$ as follows

$$V^{\pi^*} - V^{\pi_T} \leq \frac{2C_\infty^2 H^2 R_{\max}}{T} + C_\infty H \mathcal{E} + \frac{2H C_\infty |\mathcal{A}| R_{\max}}{\alpha} \varepsilon_w + \frac{4H C_\infty |\mathcal{A}| R_{\max}}{\alpha \sqrt{N}},$$

887 which completes the proof. \square

888 I.5 Proof of supporting claims

889 *Proof of Lemma H.1.* Similarly to Brukhim et al. (2022, Claim 10), the result is a direct ap-
 890 plication of Hazan and Singh (2021, Theorem 13) where we use the fact that V function is
 891 $|\mathcal{A}|R_{\max}$ -lipschitz. In Brukhim et al. (2022), the internal boosting give the following guarantee:
 892 $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d^{\pi_{t-1}}} [Q^{\pi_{t-1}}(s, \cdot)^\top \pi(\cdot | s) - Q^{\pi_{t-1}}(s, \cdot)^\top \hat{\pi}(\cdot | s)] \leq (2|\mathcal{A}|/(1-\gamma)\alpha)(\varepsilon + 2/\sqrt{N})$,
 893 where N is the number of weak learner, and $Q(\cdot, \dots)$ is bounded by $\frac{1}{1-\gamma}$ \square

894 **Lemma I.10.** *Policy gradient theorem for episodic MDPs*

$$\frac{\partial V_1^\pi}{\partial \pi_h(a_h | x_h)} = d_h^\pi(x_h) Q_h^\pi(x_h, a_h)$$

895

Proof of Lemma I.10.

$$\begin{aligned} V_1^\pi &= \mathbb{E}_{x_1 \sim \mu_1} \left[\sum_{h=1}^H r(x_h, a_h) \mid \pi \right] \\ &= \mathbb{E}_{x_1 \sim \mu_1} \left[\sum_{h'=1}^{h-1} r(x_{h'}, a_{h'}) \mid \pi_{1:h-1} \right] + \mathbb{E}_{x_h \sim d_h^\pi} \left[\sum_{h'=h}^H r(x_{h'}, a_{h'}) \mid \pi_{h:H} \right] \\ &= \mathbb{E}_{x_1 \sim \mu_1} \left[\sum_{h'=1}^{h-1} r(x_{h'}, a_{h'}) \mid \pi_{1:h-1} \right] + \mathbb{E}_{x_h \sim d^{\pi_{1:h-1}}} \left[\sum_{a \in \mathcal{A}} \pi_h(a | x_h) Q(x_h, a; \pi_{h+1:H}) \right] \end{aligned}$$

896 Since the first term has no dependence on π_h , taking derivative w.r.t. $\pi_h(a_h, x_h)$ will set it to zero.
 897 For the second term, we have d_h^π is the occupancy measure of executing policy $\pi_{1:h-1}$. Therefore,
 898 we can write the second term as

$$\mathbb{E}_{x_h \sim d_h^\pi} \left[\sum_{h'=h}^H r(x_{h'}, a_{h'}) \mid \pi_{h:H} \right] = \mathbb{E}_{x_h \sim d^{\pi_{1:h-1}}} \left[\sum_{a \in \mathcal{A}} \pi_h(a \mid x_h) Q(x_h, a; \pi_{h+1:H}) \right]$$

899 Take the partial derive w.r.t. $\pi_h(a \mid x)$ completes the proof. \square

900 *Proof of Lemma I.9.* Apply performance different lemma for $V_1(\pi^*) - V_1(\pi)$, we have

$$\begin{aligned} V_1(\pi^*) - V_1(\pi) &= \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\pi^*}} \left[\sum_{a \in \mathcal{A}} Q^\pi(x_h, a; \pi) (\pi^*(a \mid x_h) - \pi_h(a \mid x_h)) \right] \\ &= \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\pi^*}} \left[\underbrace{\sum_{a \in \mathcal{A}} Q_h^\pi(x_h, a) \pi^*(a \mid x_h) - \max_a Q_h^\pi(x_h, a)}_{< 0} \right. \\ &\quad \left. + \max_a Q^\pi(x_h, a) - \sum_{a \in \mathcal{A}} Q_h^\pi(x_h, a) \pi_h(a \mid x_h) \right] \\ &\leq \max_{h \in [H]} \max_{x_h \in \mathcal{X}_h} \frac{d_h^{\pi^*}(x_h)}{d_h^\pi(x_h)} \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^\pi} \left[\max_a Q_h^\pi(x_h, a) - \sum_{a \in \mathcal{A}} Q_h^\pi(x_h, a) \pi_h(a \mid x_h) \right], \end{aligned} \tag{25}$$

901 in the last inequality we can take the mismatch ratio $C_\infty = \max_{h \in [H]} \max_{x_h \in \mathcal{X}_h} \frac{d_h^{\pi^*}(x_h)}{d_h^\pi(x_h)}$ outside
 902 since the term inside expectation is non-negative. Additionally, for any $h \in [H]$, the expectation can
 903 be written as,

$$\begin{aligned} &\mathbb{E}_{x_h \sim d_h^\pi} \left[\max_a Q(x_h, a) - \sum_{a \in \mathcal{A}} Q(x_h, a) \pi_h(a \mid x_h) \right] \\ &\quad + \max_{\tilde{\pi}_h \in \Pi} \mathbb{E}_{x_h \sim d_h^\pi} \left[\sum_{a \in \mathcal{A}} Q(x_h, a) (\tilde{\pi}_h(a \mid x_h) - \pi_h(a \mid x_h)) \right] \\ &\leq \mathcal{E} + \mathbb{E}_{x_h \sim d_h^\pi} \left[\sum_{a \in \mathcal{A}} Q(x_h, a) (\pi'_h(a \mid x_h) - \pi(a \mid x_h)) \right] \end{aligned}$$

904 Apply Lemma I.10, we have

$$\langle \nabla_\pi V^\pi, \pi' - \pi \rangle = \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^\pi} \left[\sum_{a \in \mathcal{A}} Q_h^\pi(x_h, a) (\pi'_h(a \mid x_h) - \pi_h(a \mid x_h)) \right]$$

905 which concludes the proof. \square

906 *Proof of Lemma I.8.* Apply the performance difference Lemma L.2 and Lemma I.10, we have

$$\begin{aligned}
& \left| V^{\pi^{(1)}} - V^{\pi^{(2)}} - \langle \nabla_{\pi} V^{\pi^{(2)}}, \pi^{(1)} - \pi^{(2)} \rangle \right| \\
&= \left| \sum_{h=1}^H \mathbb{E}_{x_h \sim d^{\pi^{(1)}}_{1:h-1}} \left[\sum_{a \in \mathcal{A}} Q(x_h, a; \pi_{h+1:H}) (\pi_h^{(1)}(a | x_h) - \pi_h^{(2)}(a | x_h)) \right] \right. \\
&\quad \left. - \sum_{h=1}^H \mathbb{E}_{x_h \sim d^{\pi^{(2)}}_{1:h-1}} \left[\sum_{a \in \mathcal{A}} Q(x_h, a; \pi_{h+1:H}) (\pi_h^{(1)}(a | x_h) - \pi_h^{(2)}(a | x_h)) \right] \right| \\
&\leq \sum_{h=1}^H \left| \sum_{x_h \in \mathcal{X}_h} (d^{\pi^{(1)}}_{1:h-1}(x_h) - d^{\pi^{(2)}}_{1:h-1}(x_h)) \sum_{a \in \mathcal{A}} Q(x_h, a; \pi_{h+1:H}) (\pi_h^{(1)}(a | x_h) - \pi_h^{(2)}(a | x_h)) \right| \\
&\leq R_{\max} \sum_{h=1}^H \left| \sum_{x_h \in \mathcal{X}_h} d^{\pi^{(1)}}_{1:h-1}(x_h) - d^{\pi^{(2)}}_{1:h-1}(x_h) \right| \max_{h \in [H]} \max_{x_h \in \mathcal{X}_h} \left| \sum_{a \in \mathcal{A}} \pi_h^{(1)}(a | x_h) - \pi_h^{(2)}(a | x_h) \right| \\
&\leq R_{\max} \sum_{h=1}^H \left| \sum_{x_h \in \mathcal{X}_h} d^{\pi^{(1)}}_{1:h-1}(x_h) - d^{\pi^{(2)}}_{1:h-1}(x_h) \right| \max_{h \in [H]} \left\| \pi^{(1)} - \pi^{(2)} \right\|_{\infty, 1},
\end{aligned}$$

907 where we use the fact that $Q(x_h, a; \pi_{h+1:H}) \leq R_{\max}$. Now, for any given $h \in [H]$,

$$\begin{aligned}
& \left| \sum_{x_h \in \mathcal{X}_h} d^{\pi^{(1)}}_{1:h-1}(x_h) - d^{\pi^{(2)}}_{1:h-1}(x_h) \right| \\
&= \left| \sum_{x_h \in \mathcal{X}_h} \sum_{x_{h-1} \in \mathcal{X}_{h-1}} P(x_h | x_{h-1}, a) \pi_h(a | x_{h-1}) d^{\pi^{(1)}}_{1:h-1}(x_{h-1}) \right. \\
&\quad \left. - \sum_{x_h \in \mathcal{X}_h} \sum_{x_{h-1} \in \mathcal{X}_{h-1}} P(x_h | x_{h-1}, a) \pi_h(a | x_{h-1}) d^{\pi^{(2)}}_{1:h-1}(x_{h-1}) \right| \\
&\leq H \left\| \pi_h^{(1)} - \pi_h^{(2)} \right\|_{\infty, 1}
\end{aligned}$$

908 Therefore,

$$\sum_{h=1}^H \left| \sum_{x_h \in \mathcal{X}_h} d^{\pi^{(1)}}_{1:h-1}(x_h) - d^{\pi^{(2)}}_{1:h-1}(x_h) \right| \leq H^2 \left\| \pi^{(1)} - \pi^{(2)} \right\|_{\infty, 1},$$

909 which completes the proof. \square

910 J Discussion on Policy Cover

911 In this section, we provide a brief discussion on the policy cover with boosting. Note that boosting
912 is also a type of ensemble method, where we aggregate many policies to form a final policy. It
913 would be interesting to ask if the policy from boosting can be used as a policy cover that can well
914 explored the state space. Our first observation is that if the policies share the same partial policy in
915 the previous layer $h-1$, then we might be able to use boosting to decouple the occupancy measure
916 h ; we formalized this property in the following claim.

917 **Claim 4.** Let $\pi_{1:h-1}$ denote the partial policy over the first $h-1$ layers, and consider the boosted
918 policy at layer h defined as $\pi_h^t = \eta_t \pi_h^{t-1} + (1 - \eta_t) \hat{\pi}_h^t$. Then, for any $x_h \in \mathcal{X}_h$ the corresponding
919 occupancy measure satisfies $d_h^{\pi_{1:h-1} \circ_h \pi_h^t}(x) = (1 - \eta_t) d_h^{\pi_{1:h-1} \circ_h \pi_h^{t-1}}(x) + \eta_t d_h^{\pi_{1:h-1} \circ_h \hat{\pi}_h^t}(x)$.

920 While this property is promising and suggests that boosting can be applied in a layer-wise manner,
921 however, the resulting boosted policy remains a *single policy*, which may not be as robust as a set
922 of separate policy in a standard policy cover. Furthermore, ensemble method in boosting does not
923 guarantee the occupancy decoupling as we shown in the claim below.

924 **Claim 5.** Consider a set of non-stationary policies π^1, \dots, π^n , and let $\{\alpha_i\}_{i=1}^n \in \Delta^{n-1}$ be a set of
 925 weight on the simplex. Then, for the following ensemble method:

- 926 • **Mixture policy ensemble** (Amortila et al., 2024). Let p be the distribution over π^1, \dots, π^n
 927 given by the weights $\{\alpha_i\}_{i=1}^n$, the induced occupancy satisfies, $d^p(x) = \sum_{i \in [n]} \alpha_i d^{\pi^i}(x)$.
- 928 • **Boosting policy ensemble.** Let $\bar{\pi} = \sum_{i \in [n]} \alpha_i \pi^i$ be the boosted policy, this does not implies
 929 $d^{\bar{\pi}}(x) = \sum_{i \in [n]} \alpha_i d^{\pi^i}(x)$.

930 The fundamental difference between the two ensemble methods is that the mixture policy ensemble
 931 maintains a set of policies, and samples a policy from this set to roll out the trajectory, whereas the
 932 boosted policy ensemble is a single policy that is the convex combination of the ensemble of many
 933 policies. Although both methods are ensemble methods, at first glance, the boosting approach seems
 934 more restricted since it only outputs a single valid policy at the end. As shown by Misra et al. (2020),
 935 a single policy might fail to achieve sufficient coverage using their minimum reachability argument.
 936 This suggests that we might still need to maintain a set of policies to ensure good coverage. Our
 937 analysis arrives at the same conclusion from the perspective of boosting.

938 *Proof of Claim 4.* We write $d^\pi(x) = P(x; \pi)$ to avoid the long superscript, then for any $x_h \in \mathcal{X}$,

$$\begin{aligned}
 P(x_h; \pi_{1:h-1} \circ \pi_h^t) &= \sum_{x_{h-1}} \sum_{a \in \mathcal{A}} P(x_h | x_{h-1}, a) P(x_{h-1}; \pi_{1:h-1}) \pi_h^t(a | x_{h-1}) \\
 &= \sum_{x_{h-1}} \sum_{a \in \mathcal{A}} P(x_h | x_{h-1}, a) P(x_{h-1}; \pi_{1:h-1}) \left((1 - \eta_t) \pi_h^{t-1}(a | x_{h-1}) + \eta_t \hat{\pi}_h^t(a | x_{h-1}) \right) \\
 &= (1 - \eta_t) \sum_{x_{h-1}} \sum_{a \in \mathcal{A}} P(x_h | x_{h-1}, a) P(x_{h-1}; \pi_{1:h-1}) \pi_h^{t-1}(a | x_{h-1}) \\
 &\quad + \eta_t \sum_{x_{h-1}} \sum_{a \in \mathcal{A}} P(x_h | x_{h-1}, a) P(x_{h-1}; \pi_{1:h-1}) \hat{\pi}_h^t(a | x_{h-1}) \\
 &= (1 - \eta_t) P(x_h; \pi_{1:h-1} \circ \pi_h^{t-1}) + \eta_t P(x_h; \pi_{1:h-1} \circ \hat{\pi}_h^t),
 \end{aligned}$$

939 which completes the proof. □

940 *Proof of Claim 5.*

941 **Mixture policy ensemble.** This can be easily verified since we pick a policy w.r.t. Δ^{n-1} .

942 **Boosted (averaged) policy ensemble.** To prove this second point, we show that there exist a MDP
 943 and a set of policies $\{\pi_i\}_{i=1}^n$, that $d^{\bar{\pi}}(x) \neq \sum_{i \in [n]} \alpha_i d^{\pi_i}(x)$, where $\bar{\pi} = \sum_{i \in [n]} \alpha_i \pi_i$. Let
 944 consider an MDP with $H=3$, and two policies π_1 and π_2 . Let $\alpha \in (0, 1)$ and define the boosted
 945 policy as: $\pi = \alpha \pi_1 + (1 - \alpha) \pi_2$. The first state is sampled from the initial distribution, $x_1 \sim P_0$.
 946 For any $x_2 \in \mathcal{X}_2$, we can show that the occupancy of the boosted policy can be decouple as
 947 $d^{\bar{\pi}}(x_2) = \alpha d^{\pi_1}(x_2) + (1 - \alpha) d^{\pi_2}(x_2)$ by simply applying the definition of occupancy measure:

$$\begin{aligned}
 d^{\bar{\pi}}(x_2) &= \sum_{x_1 \in \mathcal{X}_1} P(x_2 | x_1, a) \sum_a (\alpha \pi_1(a | x_1) + (1 - \alpha) \pi_2(a | x_1)) P(x_1) \\
 &= \alpha \sum_{x_1 \in \mathcal{X}_1} P(x_2 | x_1, a) \sum_a \pi_1(a | s_1) P(x_1) + (1 - \alpha) \sum_{x_1 \in \mathcal{X}_1} P(x_2 | x_1, a) \pi_2(a | s_1) P(x_1) \\
 &= \alpha d^{\pi_1}(x_2) + (1 - \alpha) d^{\pi_2}(x_2)
 \end{aligned}$$

948 Now, for $d^{\bar{\pi}}(x_3)$,

$$\begin{aligned}
 d^{\bar{\pi}}(x_3) &= \sum_{x_2 \in \mathcal{X}_2} \sum_a P(x_3 | x_2, a) (\alpha \pi_1(a | x_2) + (1 - \alpha) \pi_2(a | x_2)) d^{\bar{\pi}}(x_2) \\
 &= \alpha^2 d^{\pi_1}(x_3) + (1 - \alpha)^2 d^{\pi_2}(x_3) \\
 &\quad + \alpha(1 - \alpha) \sum_{x_2 \in \mathcal{X}_2} P(x_3 | x_2, a) \sum_a (\pi_1(a | x_2) d^{\pi_2}(x_2) + \pi_2(a | x_2) d^{\pi_1}(x_2)).
 \end{aligned}$$

949 Let us consider the occupancy measure corresponding to $d^{\pi_1}(x_3)$:

$$\alpha^2 d^{\pi_1}(x_3) + \alpha(1 - \alpha) \sum_{x_2 \in \mathcal{X}_2} P(x_3 | x_2, a) \sum_a \pi_2(a | x_2) d^{\pi_1}(x_2),$$

950 In general, this term is not equal to $\alpha d^{\pi_1}(x_3)$ since the summation over the action depends on π_2 ;
 951 this argument also completes the proof. \square

952 We can see that even with two policies, the occupancy can be coupled and become very complicated
 953 as we move forward, with just a few more layers.

954 K Experiments

955 While the primary contribution of this work is theoretical, we include empirical evaluations to serve as
 956 a sanity check for our theory. We leave comprehensive experiments in a more complex environment
 957 for future work. We validate our algorithm on two classical control tasks from the OpenAI Gym
 958 benchmark: CartPole and MountainCar (Brockman et al., 2016). In CartPole, the agent applies a left
 959 or right force to a cart in order to balance an upright pole for as long as possible. In MountainCar,
 960 the agent must learn to first climb the left hill to build sufficient momentum to reach the goal at the top of
 961 the right hill.

962 We discretize the continuous state space of each environment into a finite number of bins along
 963 each dimension. Based on this discretization, we estimate the occupancy measure d^π using a
 964 count-based estimation over the discretized space, which is then used to construct the policy cover
 965 and the corresponding exploration bonus. We employ a depth-3 decision tree implemented using
 966 Scikit-learn (Pedregosa et al., 2011) as the weak learner. We use a total of 50 weak learners in our
 967 experiments. At each boosting iteration, the weak learner is trained to optimize the bonus-augmented
 968 objective induced by the current policy cover. We use the L_∞ coverability in the experiment, and the
 969 bonus is $b(x, a) = \frac{2}{1-\gamma} \frac{1}{\sum_{i < t} d^{\pi^i}(x, a) + 1}$ since $\tilde{C}_\infty = |\mathcal{X}| |\mathcal{A}|$, and $\mu(x, a) = \frac{1}{|\mathcal{X}| |\mathcal{A}|}$.

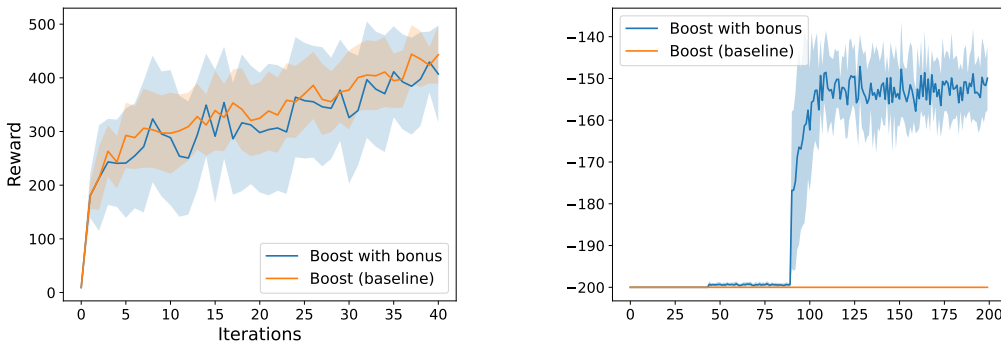


Figure 2: The total reward of cartpole (left) and mountain car (right). We compare our method with the baseline boosting of Brukhim et al. (2022). We run each experiment 10 times using boosting with 50 decision trees as weak learners, and report the total reward along with confidence intervals. Our method is competitive with the base-line on the dense-reward environment such as CartPole. For MountainCar, a sparse reward environment, our method succeeds while the baseline failed.

970 We also provide an example run of boosting MountainCar in Figure 3, along with a heatmap
 971 illustration of how the policy explores the state space. The MountainCar environment state consists
 972 of the car’s position and velocity, which we discretize into 10 bins per dimension, resulting in 100
 973 discrete states in total. As shown in the figure, the number of unique states increases steadily as
 974 more policies are added to the policy cover. The heat map highlights the regions of the state space
 975 visited by the policy, indicating that the policy gradually expands its coverage while simultaneously
 976 focusing on trajectories that lead toward the goal, balancing exploration and exploitation. All of our
 977 experiments run on a machine with 8 cores CPU and 32GB RAM.

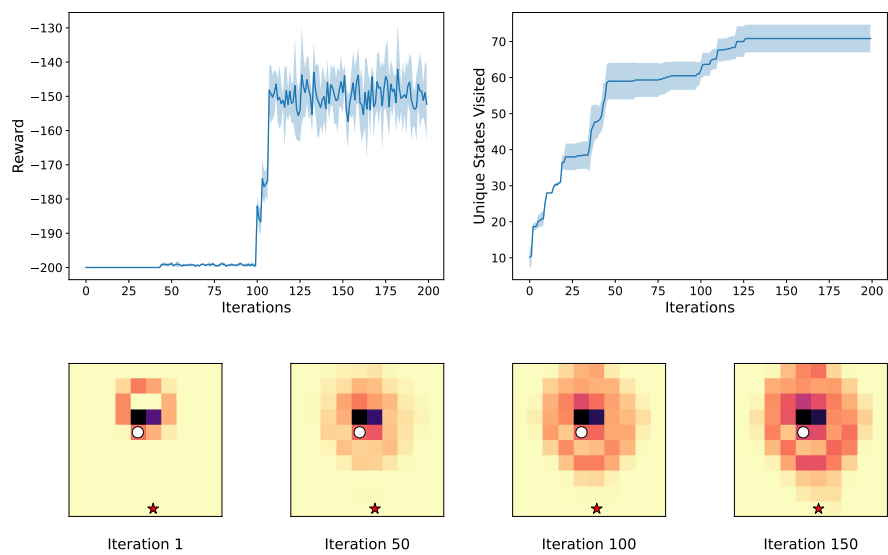


Figure 3: MountainCar results with $N = 50$ weak learners. **Top:** average episodic reward and number of unique states visited. **Bottom:** state visitation heatmaps, where the x-axis denotes velocity, and the y-axis denotes position. The start state is marked by \circ , and the goal state at zero velocity is marked by \star .

978 **L Supporting lemmas**

979 **Lemma L.1** (The performance difference lemma for discounted MDPs, (Kakade and Langford,
980 2002)). For all policies π, π' and states x_0 ,

$$V^\pi(x_0) - V^{\pi'}(x_0) = \frac{1}{1-\gamma} \mathbb{E}_{x \sim d_{x_0}^\pi} \mathbb{E}_{a \sim \pi(\cdot|x)} [A^{\pi'}(x, a)]$$

981

982 **Lemma L.2** (Performance Difference Lemma for episodic MDPs, Lemma 13 (Misra et al., 2020),).
983 For any episodic decision process with any reward function R , and any two non-stationary policies
984 $\pi_{1:H}^{(1)}$ and $\pi_{1:H}^{(2)}$, let $d_h^{\pi^{(1)}} \in \Delta(\mathcal{X}_h)$ be the distribution at time h induced by policy $\pi_{1:H}^{(1)}$. Then we
985 have

$$V(\pi_{1:H}^{(1)}) - V(\pi_{1:H}^{(2)}) = \sum_{h=1}^H \mathbb{E}_{x_h \sim d_h^{\pi^{(1)}}} \left[V(x_h; \pi_h^{(1)} \circ \pi_{h+1:H}^{(2)}) - V(x_h; \pi_{h:H}^{(2)}) \right].$$

986

987 **Lemma L.3** (Lemma 4 (Xie et al., 2023)). Let d^1, d^2, \dots, d^K be an arbitrary sequence of dis-
988 tributions over a set \mathcal{Z} , and let $\mu \in \Delta(\mathcal{Z})$ be a distribution such that $d^k(z)/\mu(z) \leq C$ for all
989 $(z, t) \in \mathcal{Z} \times [T]$. Then for all $z \in \mathcal{Z}$, we have

$$\sum_{t=1}^K \frac{d^t(z)}{\sum_{i < k} d^i(z) + C \cdot \mu(z)} \leq 2 \log(2K)$$

990 **Lemma L.4** (Lemma C.11 Amortila et al. (2024)). For any $\pi \in \Pi_{\text{rns}}, d \in \mathbb{R}_+^{\mathcal{X}}, \mu \in \mathbb{R}_+^{\mathcal{X}}$, and
991 $\varepsilon, \delta > 0$, we have that

$$\mathbb{E}^\pi \left[\frac{P_{h-1}(x_h | x_{h-1}, a_{h-1})}{d(x_h) + \varepsilon \cdot P_{h-1}(x_h | x_{h-1}, a_{h-1})} \right] \leq \mathbb{E}^\pi \left[\frac{P_{h-1}^M(x_h | x_{h-1}, a_{h-1})}{d(x_h) + \delta \cdot \mu(x_h)} \right] + \frac{\delta}{\varepsilon} \mathbb{E}^\pi \left[\frac{\mu(x_h)}{d(x_h) + \delta \cdot \mu(x_h)} \right].$$

992 **Lemma L.5** (Lemma J.2 Amortila et al. (2024)). For any distribution $\omega \in \Delta(\mathcal{Z})$ and any pair of
993 functions $w, w' : \mathcal{Z} \rightarrow \mathbb{R}_+$,

$$\mathbb{E}_\omega[w] \leq 3\mathbb{E}_\omega[w'] + 2\mathbb{E}_\omega \left[\left(\sqrt{w} - \sqrt{w'} \right)^2 \right],$$

if $w(z), w'(z) \in [0, B], \forall z \in \mathcal{Z}$ then,

$$|\mathbb{E}_\omega[w] - \mathbb{E}_\omega[w']| \leq 2B \sqrt{\mathbb{E}_\omega[(\sqrt{w} - \sqrt{w'})^2]}.$$

994

995 **Lemma L.6** (Lemma J.3 (Amortila et al., 2024)). Consider a set \mathcal{Z} and a sequence of distributions
996 $d^1, \dots, d^K \in \Delta(\mathcal{Z})$ for which there exists a distribution $\mu \in \Delta(\mathcal{Z})$ such that $\sup_{z \in \mathcal{Z}} \left\{ \frac{d^k(z)}{\mu(z)} \right\} \leq C$
997 for all $k \in [K]$. For any sequence of functions $g^1, \dots, g^K \subset (\mathcal{Z} \rightarrow [-B, B])$, it holds that

$$\sum_{k=1}^K \mathbb{E}_{z \sim d^k} [g(z)] \leq \sqrt{2C \log(2K) \sum_{k=1}^K \sum_{i < k} \mathbb{E}_{z \sim d^i} \left[(g^i(z))^2 \right]} + 2CB.$$

998 **L.1 Boosting Framework (Brukhim et al., 2022)**

999 In this section, we provide the subroutines from Brukhim et al. (2022) that were omitted from the
1000 main text. Internal Boost (Algorithm 12) instantiates Algorithm 3 from Hazan and Singh (2021). At
1001 a high level, Internal Boost can be viewed as an *offline stochastic contextual optimization* problem;
1002 it is a straightforward application of Theorem 13 from Hazan and Singh (2021), and we restate its
1003 guarantees in the following lemma:

1004 **Lemma L.7** (Claim 10 Brukhim et al. (2022)). *Let $\beta = \sqrt{1/\alpha N}$, $\eta_{2,n} = \min\{2/n, 1\}$, $\|Q\| \in$
1005 $[0, B]$. π'_t produced by Algorithm 1 satisfies $\max_{\pi \in \Pi} \mathbb{E}_{(s,Q) \sim \mathcal{D}_t} [Q^\top \pi(s)] -$
1006 $\mathbb{E}_{(s,Q) \sim \mathcal{D}_t} [Q^\top \pi'_t(s)] \leq \frac{2B}{\alpha} (\varepsilon_w + 2/\sqrt{N})$.*

Algorithm 12 Internal Boost (Algorithm 2 of Brukhim et al. (2022))

- 1: **Input:** number of iterations N , number of episodes M , initial policy π , initial state distribution μ .
- 2: Set $\tilde{\pi}_0$ to be an arbitrary policy in Π .
- 3: **for** $n = 1$ **to** N **do**
- 4: Execute π with μ via Algorithm 13 for M episodes, to get $D_n = \{(s_i, \widehat{Q}_i)_{i=1}^M\}$.
- 5: Modify D_n to produce a new dataset $D'_n = \{(s_i, f_i)\}_{i=1}^M$, such that for all $i \in [m]$:

$$f_i = \frac{1}{\beta} (y_i - \tilde{\pi}_n(\cdot|s_i)),$$

$$y_i = \operatorname{argmin}_{y \in \mathbb{R}^{|\mathcal{A}|}} \left\{ -\widehat{Q}_i^\top y + G \min_{z \in \Delta_{\mathcal{A}}} \|z - y\| + \frac{\|\tilde{\pi}_n(\cdot|s_i) - y\|^2}{2\beta} \right\},$$

where $G = \frac{|\mathcal{A}|}{1-\gamma}$, $\beta = \frac{2\gamma}{(1-\gamma)^3}$ and $f_i, \widehat{Q}_i \in \mathbb{R}^{|\mathcal{A}|}$.

- 6: Let \mathcal{L}_n be the policy chosen by the weak learning oracle when given data set D'_n .
- 7: Update

$$\tilde{\pi}_n = (1 - \eta_{2,n})\tilde{\pi}_{n-1} + \frac{\eta_{2,n}}{\alpha} \mathcal{L}_n.$$

- 8: **end for**
 - 9: **return** $\Gamma[\tilde{\pi}_N]$.
-

Algorithm 13 Trajectory Sampler: samples a state $s \sim d^\pi$, and an unbiased estimate of Q_s^π

- 1: Sample state $s_0 \sim \mu$, action $a' \sim \mathcal{U}(\mathcal{A})$ uniformly.
 - 2: Sample $s \sim d^\pi$ as follows: at every timestep h , with probability γ , act according to π ; else, accept s_h as the sample and proceed to Step 3.
 - 3: Take action a' at state s_h , then continue to execute π , and use a termination probability of $1 - \gamma$. Upon termination, set $R(s_h, a')$ as the *undiscounted* sum of rewards from time h onwards.
 - 4: Define the vector $\widehat{Q}_{s_h}^\pi$, such that for all $a \in \mathcal{A}$, $\widehat{Q}_{s_h}^\pi(a) = |\mathcal{A}| \cdot R(s_h, a') \cdot \mathbb{1}_{a=a'}$.
 - 5: **return** $(s_h, \widehat{Q}_{s_h}^\pi)$.
-

1007 An important generalization of the property of convexity is gradient domination, which we restate
1008 below.

Definition L.8 (Gradient Domination). A function $f : \mathcal{K} \rightarrow \mathbb{R}$ is said to be $(\kappa, \tau, \mathcal{K}_1, \mathcal{K}_2)$ -locally gradient dominated (around \mathcal{K}_1 by \mathcal{K}_2) if for all $x \in \mathcal{K}_1$, it holds that

$$\max_{y \in \mathcal{K}} f(y) - f(x) \leq \kappa \cdot \max_{y \in \mathcal{K}_2} \left\{ \nabla f(x)^\top (y - x) \right\} + \tau.$$

1009 The following smoothness property of the V function is also helpful for our analysis that we use on
1010 Appendix D and Appendix E.

1011 **Lemma L.9** (Lemma 16 (Brukhim et al., 2022)). V^π is $\frac{2B\gamma}{(1-\gamma)^2}$ -smooth in the $\|\cdot\|_{\infty,1}$ norm, for
1012 $V^\pi(x) \in [0, B]$, $\forall x \in \mathcal{X}$.

1013 The boosting in Algorithm 1 can be cast as a non-convex Frank-Wolfe algorithm, which we provide
 1014 the pseudo-code in Algorithm 14, we denote by \mathcal{O} a black-box oracle to an $(\epsilon_0, \mathcal{K}_2)$ -approximate
 1015 linear optimizer over a convex set $\mathcal{K} \subseteq \mathbb{R}^d$ such that for any given $v \in \mathbb{R}^d$, we have $v^\top \mathcal{O}(v) \geq$
 1016 $\max_{u \in \mathcal{K}_2} v^\top u - \epsilon_0$.

Algorithm 14 Non-convex Frank-Wolfe

- 1: Input: $T > 0$, objective f , linear optimizer \mathcal{O} , rate η_t .
 - 2: Choose $x_0 \in \mathcal{K}$ arbitrarily.
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Call $z_t = \mathcal{O}(\nabla_{t-1})$, where $\nabla_{t-1} = \nabla f(x_{t-1})$. Set $x_t = (1 - \eta_t)x_{t-1} + \eta_t z_t$.
 - 5: **end for**
 - 6: **return** $\bar{x} := x_{t'}$ where $t' = \operatorname{argmin}_t \nabla_t^\top (z_t - x_t)$.
-

1017 We provide the guarantee of Algorithm 14 as follows:

Lemma L.10 (Theorem 9 (Bruckhim et al., 2022)). *Let $f : \mathcal{K} \rightarrow \mathbb{R}$ be L -smooth in some norm $\|\cdot\|_*$, bounded for all $x \in \mathcal{K}$, $|f(x)| \leq H$ for some $H > 0$, and let the diameter of \mathcal{K} in $\|\cdot\|_*$ be D . Then, for a $(\epsilon_0, \mathcal{K}_2)$ -linear optimization oracle \mathcal{O} , and $\eta_t = \eta = \sqrt{\frac{4H}{LD^2T}}$, the output \bar{x} of Algorithm 14 satisfies*

$$\max_{x^* \in \mathcal{K}} f(x^*) - f(\bar{x}) \leq \frac{2\kappa^2 \max\{LD^2, H\}}{T} + \tau + \kappa\epsilon_0$$

1018 Furthermore, if f is $(\kappa, \tau, \mathcal{K}_1, \mathcal{K}_2)$ -locally gradient-dominated and $x_0, \dots, x_T \in \mathcal{K}_1$, then the output
 1019 \bar{x} of Algorithm 14 where $\eta_t = \min\{1, \frac{2\kappa}{t}\}$ satisfies the bound on the right.

1020 **NeurIPS Paper Checklist**

1021 **1. Claims**

1022 Question: Do the main claims made in the abstract and introduction accurately reflect the
1023 paper’s contributions and scope?

1024 Answer: **[Yes]**

1025 Justification: The main theoretical results in Section 3 directly support the claims stated in
1026 the abstract and introduction. In particular, we provide guarantees that relax the bounded
1027 distribution mismatch assumption required in prior RL boosting work, which constitutes the
1028 central contribution of the paper.

1029 Guidelines:

- 1030 • The answer **[N/A]** means that the abstract and introduction do not include the claims
1031 made in the paper.
- 1032 • The abstract and/or introduction should clearly state the claims made, including the
1033 contributions made in the paper and important assumptions and limitations. A **[No]** or
1034 **[N/A]** answer to this question will not be perceived well by the reviewers.
- 1035 • The claims made should match theoretical and experimental results, and reflect how
1036 much the results can be expected to generalize to other settings.
- 1037 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1038 are not attained by the paper.

1039 **2. Limitations**

1040 Question: Does the paper discuss the limitations of the work performed by the authors?

1041 Answer: **[Yes]**

1042 Justification: We discuss limitations and possibilities for future work in the second paragraph
1043 of the conclusion section.

1044 Guidelines:

- 1045 • The answer **[N/A]** means that the paper has no limitation while the answer **[No]** means
1046 that the paper has limitations, but those are not discussed in the paper.
- 1047 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 1048 • The paper should point out any strong assumptions and how robust the results are to
1049 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1050 model well-specification, asymptotic approximations only holding locally). The authors
1051 should reflect on how these assumptions might be violated in practice and what the
1052 implications would be.
- 1053 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1054 only tested on a few datasets or with a few runs. In general, empirical results often
1055 depend on implicit assumptions, which should be articulated.
- 1056 • The authors should reflect on the factors that influence the performance of the approach.
1057 For example, a facial recognition algorithm may perform poorly when image resolution
1058 is low or images are taken in low lighting. Or a speech-to-text system might not be
1059 used reliably to provide closed captions for online lectures because it fails to handle
1060 technical jargon.
- 1061 • The authors should discuss the computational efficiency of the proposed algorithms
1062 and how they scale with dataset size.
- 1063 • If applicable, the authors should discuss possible limitations of their approach to
1064 address problems of privacy and fairness.
- 1065 • While the authors might fear that complete honesty about limitations might be used by
1066 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1067 limitations that aren’t acknowledged in the paper. The authors should use their best
1068 judgment and recognize that individual actions in favor of transparency play an impor-
1069 tant role in developing norms that preserve the integrity of the community. Reviewers
1070 will be specifically instructed to not penalize honesty concerning limitations.

1071 **3. Theory assumptions and proofs**

1072 Question: For each theoretical result, does the paper provide the full set of assumptions and
1073 a complete (and correct) proof?

1074 Answer: [Yes]

1075 Justification: We explain the detailed setup in Section 2, provide all the assumptions along
1076 with the main results in Theorem 3.2 and Theorem 3.5, we provide complete proofs for each
1077 result in the Appendix section.

1078 Guidelines:

- 1079 • The answer [N/A] means that the paper does not include theoretical results.
- 1080 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
1081 referenced.
- 1082 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1083 • The proofs can either appear in the main paper or the supplemental material, but if
1084 they appear in the supplemental material, the authors are encouraged to provide a short
1085 proof sketch to provide intuition.
- 1086 • Inversely, any informal proof provided in the core of the paper should be complemented
1087 by formal proofs provided in appendix or supplemental material.
- 1088 • Theorems and Lemmas that the proof relies upon should be properly referenced.

1089 4. Experimental result reproducibility

1090 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
1091 perimental results of the paper to the extent that it affects the main claims and/or conclusions
1092 of the paper (regardless of whether the code and data are provided or not)?

1093 Answer: [Yes]

1094 Justification: The experiment section provides the necessary details for replication, including
1095 the specific type of the weak learner and the construction of the exploration bonus.

1096 Guidelines:

- 1097 • The answer [N/A] means that the paper does not include experiments.
- 1098 • If the paper includes experiments, a [No] answer to this question will not be perceived
1099 well by the reviewers: Making the paper reproducible is important, regardless of
1100 whether the code and data are provided or not.
- 1101 • If the contribution is a dataset and/or model, the authors should describe the steps taken
1102 to make their results reproducible or verifiable.
- 1103 • Depending on the contribution, reproducibility can be accomplished in various ways.
1104 For example, if the contribution is a novel architecture, describing the architecture fully
1105 might suffice, or if the contribution is a specific model and empirical evaluation, it may
1106 be necessary to either make it possible for others to replicate the model with the same
1107 dataset, or provide access to the model. In general, releasing code and data is often
1108 one good way to accomplish this, but reproducibility can also be provided via detailed
1109 instructions for how to replicate the results, access to a hosted model (e.g., in the case
1110 of a large language model), releasing of a model checkpoint, or other means that are
1111 appropriate to the research performed.
- 1112 • While NeurIPS does not require releasing code, the conference does require all submis-
1113 sions to provide some reasonable avenue for reproducibility, which may depend on the
1114 nature of the contribution. For example
 - 1115 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
1116 to reproduce that algorithm.
 - 1117 (b) If the contribution is primarily a new model architecture, the paper should describe
1118 the architecture clearly and fully.
 - 1119 (c) If the contribution is a new model (e.g., a large language model), then there should
1120 either be a way to access this model for reproducing the results or a way to reproduce
1121 the model (e.g., with an open-source dataset or instructions for how to construct
1122 the dataset).
 - 1123 (d) We recognize that reproducibility may be tricky in some cases, in which case
1124 authors are welcome to describe the particular way they provide for reproducibility.
1125 In the case of closed-source models, it may be that access to the model is limited in

1126 some way (e.g., to registered users), but it should be possible for other researchers
1127 to have some path to reproducing or verifying the results.

1128 5. Open access to data and code

1129 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1130 tions to faithfully reproduce the main experimental results, as described in supplemental
1131 material?

1132 Answer: [Yes]

1133 Justification: We include the script for the experiment in the supplemental material.

1134 Guidelines:

- 1135 • The answer [N/A] means that paper does not include experiments requiring code.
- 1136 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
1137 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1138 • While we encourage the release of code and data, we understand that this might not
1139 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
1140 including code, unless this is central to the contribution (e.g., for a new open-source
1141 benchmark).
- 1142 • The instructions should contain the exact command and environment needed to run to
1143 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1144 • The authors should provide instructions on data access and preparation, including how
1145 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1146 • The authors should provide scripts to reproduce all experimental results for the new
1147 proposed method and baselines. If only a subset of experiments are reproducible, they
1148 should state which ones are omitted from the script and why.
- 1149 • At submission time, to preserve anonymity, the authors should release anonymized
1150 versions (if applicable).
- 1151 • Providing as much information as possible in supplemental material (appended to the
1152 paper) is recommended, but including URLs to data and code is permitted.

1154 6. Experimental setting/details

1155 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
1156 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1157 Answer: [Yes]

1158 Justification: We include the detail information in Appendix K.

1159 Guidelines:

- 1160 • The answer [N/A] means that the paper does not include experiments.
- 1161 • The experimental setting should be presented in the core of the paper to a level of detail
1162 that is necessary to appreciate the results and make sense of them.
- 1163 • The full details can be provided either with the code, in appendix, or as supplemental
1164 material.

1165 7. Experiment statistical significance

1166 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1167 information about the statistical significance of the experiments?

1168 Answer: [Yes]

1169 Justification: We report error bars for our method as well as the baseline from prior work to
1170 ensure statistically meaningful comparison.

1171 Guidelines:

- 1172 • The answer [N/A] means that the paper does not include experiments.
- 1173 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
1174 intervals, or statistical significance tests, at least for the experiments that support the
1175 main claims of the paper.

- 1176 • The factors of variability that the error bars are capturing should be clearly stated (for
1177 example, train/test split, initialization, random drawing of some parameter, or overall
1178 run with given experimental conditions).
- 1179 • The method for calculating the error bars should be explained (closed form formula,
1180 call to a library function, bootstrap, etc.)
- 1181 • The assumptions made should be given (e.g., Normally distributed errors).
- 1182 • It should be clear whether the error bar is the standard deviation or the standard error
1183 of the mean.
- 1184 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1185 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1186 of Normality of errors is not verified.
- 1187 • For asymmetric distributions, the authors should be careful not to show in tables or
1188 figures symmetric error bars that would yield results that are out of range (e.g., negative
1189 error rates).
- 1190 • If error bars are reported in tables or plots, the authors should explain in the text how
1191 they were calculated and reference the corresponding figures or tables in the text.

1192 8. Experiments compute resources

1193 Question: For each experiment, does the paper provide sufficient information on the com-
1194 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1195 the experiments?

1196 Answer: [Yes]

1197 Justification: We report the computation resource in the Appendix K.

1198 Guidelines:

- 1199 • The answer [N/A] means that the paper does not include experiments.
- 1200 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1201 or cloud provider, including relevant memory and storage.
- 1202 • The paper should provide the amount of compute required for each of the individual
1203 experimental runs as well as estimate the total compute.
- 1204 • The paper should disclose whether the full research project required more compute
1205 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1206 didn't make it into the paper).

1207 9. Code of ethics

1208 Question: Does the research conducted in the paper conform, in every respect, with the
1209 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1210 Answer: [Yes]

1211 Justification: We have reviewed the Code of Ethics and have made every effort to adhere to
1212 it.

1213 Guidelines:

- 1214 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
1215 Ethics.
- 1216 • If the authors answer [No], they should explain the special circumstances that require a
1217 deviation from the Code of Ethics.
- 1218 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1219 eration due to laws or regulations in their jurisdiction).

1220 10. Broader impacts

1221 Question: Does the paper discuss both potential positive societal impacts and negative
1222 societal impacts of the work performed?

1223 Answer: [N/A]

1224 Justification: As this is a theoretical paper in nature, the societal impacts of this paper are
1225 somewhat limited in scope.

1226 Guidelines:

- 1227 • The answer [N/A] means that there is no societal impact of the work performed.
- 1228 • If the authors answer [N/A] or [No], they should explain why their work has no societal
- 1229 impact or why the paper does not address societal impact.
- 1230 • Examples of negative societal impacts include potential malicious or unintended uses
- 1231 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 1232 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 1233 groups), privacy considerations, and security considerations.
- 1234 • The conference expects that many papers will be foundational research and not tied
- 1235 to particular applications, let alone deployments. However, if there is a direct path to
- 1236 any negative applications, the authors should point it out. For example, it is legitimate
- 1237 to point out that an improvement in the quality of generative models could be used to
- 1238 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
- 1239 that a generic algorithm for optimizing neural networks could enable people to train
- 1240 models that generate Deepfakes faster.
- 1241 • The authors should consider possible harms that could arise when the technology is
- 1242 being used as intended and functioning correctly, harms that could arise when the
- 1243 technology is being used as intended but gives incorrect results, and harms following
- 1244 from (intentional or unintentional) misuse of the technology.
- 1245 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 1246 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 1247 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 1248 feedback over time, improving the efficiency and accessibility of ML).

1249 11. Safeguards

1250 Question: Does the paper describe safeguards that have been put in place for responsible
1251 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1252 image generators, or scraped datasets)?

1253 Answer: [N/A]

1254 Justification: As this is a mainly theoretical work in nature, it does not present the high-risk
1255 profiles described in the guidelines.

1256 Guidelines:

- 1257 • The answer [N/A] means that the paper poses no such risks.
- 1258 • Released models that have a high risk for misuse or dual-use should be released with
- 1259 necessary safeguards to allow for controlled use of the model, for example by requiring
- 1260 that users adhere to usage guidelines or restrictions to access the model or implementing
- 1261 safety filters.
- 1262 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 1263 should describe how they avoided releasing unsafe images.
- 1264 • We recognize that providing effective safeguards is challenging, and many papers do
- 1265 not require this, but we encourage authors to take this into account and make a best
- 1266 faith effort.

1267 12. Licenses for existing assets

1268 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1269 the paper, properly credited and are the license and terms of use explicitly mentioned and
1270 properly respected?

1271 Answer: [Yes]

1272 Justification: We use the OpenAI Gym in our experiments, which is cited in the experiment
1273 section, and we have complied with all terms of use.

1274 Guidelines:

- 1275 • The answer [N/A] means that the paper does not use existing assets.
- 1276 • The authors should cite the original paper that produced the code package or dataset.
- 1277 • The authors should state which version of the asset is used and, if possible, include a
- 1278 URL.
- 1279 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 1280 • For scraped data from a particular source (e.g., website), the copyright and terms of
1281 service of that source should be provided.
- 1282 • If assets are released, the license, copyright information, and terms of use in the
1283 package should be provided. For popular datasets, paperswithcode.com/datasets
1284 has curated licenses for some datasets. Their licensing guide can help determine the
1285 license of a dataset.
- 1286 • For existing datasets that are re-packaged, both the original license and the license of
1287 the derived asset (if it has changed) should be provided.
- 1288 • If this information is not available online, the authors are encouraged to reach out to
1289 the asset's creators.

1290 13. New assets

1291 Question: Are new assets introduced in the paper well documented and is the documentation
1292 provided alongside the assets?

1293 Answer: [N/A]

1294 Justification: We do not have any new assets in this work.

1295 Guidelines:

- 1296 • The answer [N/A] means that the paper does not release new assets.
- 1297 • Researchers should communicate the details of the dataset/code/model as part of their
1298 submissions via structured templates. This includes details about training, license,
1299 limitations, etc.
- 1300 • The paper should discuss whether and how consent was obtained from people whose
1301 asset is used.
- 1302 • At submission time, remember to anonymize your assets (if applicable). You can either
1303 create an anonymized URL or include an anonymized zip file.

1304 14. Crowdsourcing and research with human subjects

1305 Question: For crowdsourcing experiments and research with human subjects, does the paper
1306 include the full text of instructions given to participants and screenshots, if applicable, as
1307 well as details about compensation (if any)?

1308 Answer: [N/A]

1309 Justification: This work does not involve crowdsourcing or research with human subjects.

1310 Guidelines:

- 1311 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1312 with human subjects.
- 1313 • Including this information in the supplemental material is fine, but if the main contribu-
1314 tion of the paper involves human subjects, then as much detail as possible should be
1315 included in the main paper.
- 1316 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1317 or other labor should be paid at least the minimum wage in the country of the data
1318 collector.

1319 15. Institutional review board (IRB) approvals or equivalent for research with human 1320 subjects

1321 Question: Does the paper describe potential risks incurred by study participants, whether
1322 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1323 approvals (or an equivalent approval/review based on the requirements of your country or
1324 institution) were obtained?

1325 Answer: [N/A]

1326 Justification: This work does not involve crowdsourcing or research with human subjects.

1327 Guidelines:

- 1328 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1329 with human subjects.

- 1330 • Depending on the country in which research is conducted, IRB approval (or equivalent)
- 1331 may be required for any human subjects research. If you obtained IRB approval, you
- 1332 should clearly state this in the paper.
- 1333 • We recognize that the procedures for this may vary significantly between institutions
- 1334 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 1335 guidelines for their institution.
- 1336 • For initial submissions, do not include any information that would break anonymity (if
- 1337 applicable), such as the institution conducting the review.

1338 **16. Declaration of LLM usage**

1339 Question: Does the paper describe the usage of LLMs if it is an important, original, or

1340 non-standard component of the core methods in this research? Note that if the LLM is used

1341 only for writing, editing, or formatting purposes and does *not* impact the core methodology,

1342 scientific rigor, or originality of the research, declaration is not required.

1343 Answer: [N/A]

1344 Justification: We only use LLM for editing assistance.

1345 Guidelines:

- 1346 • The answer [N/A] means that the core method development in this research does not
- 1347 involve LLMs as any important, original, or non-standard components.
- 1348 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
- 1349 be described.