

1 In Appendix A we introduce some basic definitions that are needed for our theoretical results. In
 2 Appendix B, we provide sufficient conditions for Assumption 1 that were mentioned in the main
 3 text. In Appendix C and Appendix D we prove the error bounds for PPI and PQI. In Appendix E
 4 and Appendix F we present more details of our experimental results.

5 A Definition of auxiliary MDP and policy projection

6 First we introduce the definition of an auxiliary MDP M' based on M : each state in M has an
 7 absorbing action which leads to a self-looping absorbing state. All the other dynamics are preserved.
 8 Rewards are 0 for the absorbing action and unchanged elsewhere. More formally: The auxiliary
 9 MDP M' given $M = \langle \mathcal{S}, \mathcal{A}, R, P, \gamma, \rho \rangle$ is defined as $M' = \langle \mathcal{S}', \mathcal{A}', R', P', \gamma, \rho \rangle$, where
 10 $\mathcal{S}' = \mathcal{S} \cup \{s_{\text{abs}}\}$, $\mathcal{A}' = \mathcal{A} \cup \{a_{\text{abs}}\}$. R' and P' are the same as R and P for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.
 11 $R'(s, a)$ if $s = s_{\text{abs}}$ or $a = a_{\text{abs}}$ is a point mass on 0, and $P'(s, a)$ if $s = s_{\text{abs}}$ or $a = a_{\text{abs}}$ is a point
 12 mass on s_{abs} . A data set D generated from distribution μ on M is also from the distribution μ on
 13 M' , since all distributions on $\mathcal{S} \times \mathcal{A}$ are the same between the two MDPs. This MDP is used only to
 14 perform our analysis about the error bounds on the algorithm, and is not needed at all for executing
 15 Algorithm 1 and 2. As some of the notations is actually a function of the MDP, we clarify the usage
 16 of notation w.r.t. M/M' in the appendix:

- 17 1. Policy value functions V^π/Q^π and Bellman operators $\mathcal{T}/\mathcal{T}^\pi$ correspond to M' unless they
 18 have additional subscripts.
- 19 2. The definition of \mathcal{F} , Π , \mathcal{T}_ζ , \mathcal{T}_ζ^π , $\hat{\mu}$ is independent of the change from M to M' .
- 20 3. μ is also a distribution over $\mathcal{S}' \times \mathcal{A}'$. The definition of ζ will be extended to $\mathcal{S}' \times \mathcal{A}'$ as
 21 follow:

$$\zeta(s, a) = \begin{cases} \mathbb{1}(\hat{\mu}(s, a) \geq b) & s \in \mathcal{S}, a \in \mathcal{A} \\ 0 & s = s_{\text{abs}} \text{ or } a = a_{\text{abs}} \end{cases}$$

22 (That means there is only one version of μ and ζ across M and M' , instead of like we have
 23 $\mathcal{T}_{M'}^\pi$ and \mathcal{T}_M^π for M and M' .)

24 Recall the definition of semi-norm of any function of state-action pairs. For any function $g : \mathcal{S}' \times$
 25 $\mathcal{A}' \rightarrow \mathbb{R}$, $\nu \in \Delta(\mathcal{S}' \times \mathcal{A}')$, and $p \geq 1$, define the shorthand $\|g\|_{p, \nu} := (\mathbb{E}_{(s, a) \sim \nu}[|g(s, a)|^p])^{1/p}$. With
 26 some abuse of notation, later we also use this norm for $\nu \in \Delta(\mathcal{S} \times \mathcal{A})$ (specifically, μ) by viewing
 27 the probability of ν on additional (s, a) pairs as zero. Given a policy π , let $\eta_h^\pi(s)$ be the marginal
 28 distribution of s_h under π , that is, $\eta_h^\pi(s) := \Pr[s_h = s | s_0 \sim p, \pi]$, $\eta_h^\pi(s, a) = \eta_h^\pi(s)\pi(a|s)$,
 29 and $\eta^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \eta_h^\pi(s, a)$. We also use $P(s, a)$ and $P(\nu)$ to denote the next state
 30 distribution given a state action pair or given the current state action distribution.

31 The norm $\|\cdot\|_{p, \nu}$ are defined over $\mathcal{S}' \times \mathcal{A}'$. Though for the input space of function $f \in \mathcal{F}$ is $\mathcal{S} \times \mathcal{A}$,
 32 the norm can still be well-defined. All of the norm would not need the value of $f(s, a)$ on $s = s_{\text{abs}}$
 33 or $a = a_{\text{abs}}$, because the distribution does not cover those (s, a) , or the f inside of the norm is
 34 multiplied by other function that is zero for those (s, a) .

35 We first formally state an obvious result about policy value in M and M' .

36 **Lemma 1.** *For any policy π that only have non-zero probability for $a \in \mathcal{A}$, $v_{M'}^\pi = v_M^\pi$.*

Proof. By the definition of M' , P and R are the same with M over $\mathcal{S} \times \mathcal{A}$.

$$v_M^\pi = \mathbb{E}_M \left[\sum_{t=0}^h \gamma^t r_t | s_0 \sim p, \pi \right] = \mathbb{E}_{M'} \left[\sum_{t=0}^h \gamma^t r_t | s_0 \sim p, \pi \right] = v_{M'}^\pi$$

37 □

38 For the readability we repeat the Definition 1 here

39 **Definition 1** (ζ -constrained policy set). *Let Π_C^{all} be the set of policies $\mathcal{S} \rightarrow \Delta(\mathcal{A})$ such that*
 40 $\Pr(\zeta(s, a) = 0 | \pi) \leq \epsilon_\zeta$. *That is*

$$(1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s, a \sim \eta_h^\pi} [\mathbb{1}(\zeta(s, a) = 0)] \leq \epsilon_\zeta \quad (4)$$

41 Now we introduce another constrained policy set. Different from ζ -constrained policy set which
 42 we introduced in Definition 1, this policy set is on M' instead of M and the policy is forced to take
 43 action a_{abs} when $\zeta(s, a) = 0$ for all a . The reason we introduce this is to help us formally analyze the
 44 (lower bound of) performance of the resulting policy. We essentially treat any action taken outside
 45 of the support to be a_{abs} . Later we will define a projection to achieve that and show results about
 46 how the policy value changes after projection.

47 **Definition 2** (strong ζ -constrained policy set). *Let Π_{SC}^{all} be the set of all policies $S' \rightarrow \Delta(\mathcal{A}')$ such
 48 that for $\forall (s, a)$ $\pi(a|s) > 0$ then 1) $\zeta(s, a) > 0$, or 2) $a = a_{\text{abs}}$.*

49 Notice that for ζ -constrained policy set we have no requirement for π if for any action $\zeta(s, a)$ is zero.
 50 For strong ζ -constrained policy set we enforce π to take action a_{abs} . The second difference is ζ -
 51 constrained policy set requires the condition holds for s, a that is reachable, which means $\eta_h^\pi(s) > 0$
 52 and $\pi(a|s) > 0$. Here we require the same condition holds for any s, a such that $\pi(a|s) > 0$. In
 53 general, this is a stronger definition. However, we can show that for any policy in ζ -constrained
 54 policy set, it can be mapped to a policy in strong ζ -constrained policy set, with changing value
 55 bounds. Since we only need to change the behavior of policy in the state actions such that the state
 56 actions that $\zeta = 0$, the value of policy will not be much different.

57 Now we define a projection that maps any policy to Π_{SC}^{all} .

58 **Definition 3** (ζ -constrained policy projection). *$(\Xi\pi)(a|s)$ equals $\zeta(s, a)\pi(a|s)$ if $a \in \mathcal{A}$, and equals
 59 $\sum_{a' \in \mathcal{A}'} \pi(a'|s)(1 - \zeta(s, a'))$ if $a = a_{\text{abs}}$*

60 Next we show that the projection of policy will has an equal or smaller value than the original policy.

61 **Lemma 2.** *For any policy $\pi : S' \rightarrow \Delta(\mathcal{A}')$, $v_{M'}^\pi \geq v_{M'}^{\Xi(\pi)}$, and $v_{M'}^\pi = v_{M'}^{\Xi(\pi)}$ if for any (s, a)
 62 reachable by π , $\zeta(s, a) = 1$.*

63 *Proof.* We drop the subscription of M' in this proof for ease of notation. For any given s ,

$$\sum_{a \in \mathcal{A}'} \pi(a|s) Q^{\Xi(\pi)}(s, a) = \sum_{a \in \mathcal{A}} \pi(a|s) Q^{\Xi(\pi)}(s, a) \quad (Q^\pi(s, a_{\text{abs}} = 0))$$

$$\geq \sum_{a \in \mathcal{A}} \zeta(a|s) \pi(a|s) Q^{\Xi(\pi)}(s, a) \quad (5)$$

$$= \Xi(\pi)(a_{\text{abs}}|s) Q^{\Xi(\pi)}(s, a_{\text{abs}}) + \sum_{a \in \mathcal{A}} \Xi(\pi)(a|s) Q^{\Xi(\pi)}(s, a) \quad (\text{Def of } \Xi)$$

$$= \sum_{a \in \mathcal{A}'} \Xi(\pi)(a|s) Q^{\Xi(\pi)}(s, a) \quad (6)$$

$$= V^{\Xi(\pi)}(s) \quad (7)$$

64 The inequality is an equality if for any a s.t. $\pi(a|s) > 0$, $\zeta(s, a) = 1$. By the performance difference
 65 lemma [3, Lemma 6.1]:

$$v^{\Xi(\pi)} - v^\pi = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^\pi} \left[V^{\Xi(\pi)}(s) - \sum_{a \in \mathcal{A}'} \pi(a|s) Q^{\Xi(\pi)}(s, a) \right] \leq 0 \quad (8)$$

66 The inequality is an equality if for any (s, a) s.t. $\eta_h^\pi(s) \pi(a|s) > 0$ for some h , $\zeta(s, a) = 1$.
 67 In another word for any state-action reachable by π ($\eta_h^\pi(s) > 0$ and $\pi(a|s) > 0$ for some h),
 68 $\zeta(s, a) = 1$. \square

69 The following results shows for any policy π in the ζ -constrained policy set the projection will not
 70 change the policy value much.

71 **Lemma 3.** *For any policy $\pi \in \Pi_C^{\text{all}}$, $v_M^\pi \leq v_M^{\Xi(\pi)} + \frac{\epsilon_\zeta V_{\text{max}}}{1-\gamma}$*

72 *Proof.* Since π only takes action in \mathcal{A} , by Lemma 1, we have that $v_M^\pi = v_{M'}^\pi$. Since $\pi \in \Pi_C^{\text{all}}$, we
 73 have that $\Pr(\zeta(s, a) = 0 | \pi) \leq \epsilon_\zeta$, which means that:

$$(1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^\pi} [\mathbf{1}(\zeta(s, a) = 0)] \leq \epsilon_\zeta \quad (9)$$

74 Thus:

$$v^{\Xi(\pi)} - v^\pi = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^\pi} \left[V^{\Xi(\pi)}(s) - \sum_{a \in \mathcal{A}'} \pi(a|s) Q^{\Xi(\pi)}(s, a) \right] \quad (10)$$

$$= \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^\pi} \left[V^{\Xi(\pi)}(s) - \sum_{a \in \mathcal{A}'} \pi(a|s) \zeta(s, a) Q^{\Xi(\pi)}(s, a) \right] \quad (11)$$

$$- \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s, a \sim \eta_h^\pi} \left[\mathbf{1}(\zeta(s, a) = 0) Q^{\Xi(\pi)}(s, a) \right] \quad (12)$$

$$\geq \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^\pi} \left[V^{\Xi(\pi)}(s) - \sum_{a \in \mathcal{A}'} \pi(a|s) \zeta(s, a) Q^{\Xi(\pi)}(s, a) \right] \quad (13)$$

$$- V_{\max} \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s, a \sim \eta_h^\pi} [\mathbf{1}(\zeta(s, a) = 0)] \quad (14)$$

$$\geq \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^\pi} \left[V^{\Xi(\pi)}(s) - \sum_{a \in \mathcal{A}'} \pi(a|s) \zeta(s, a) Q^{\Xi(\pi)}(s, a) \right] - \frac{V_{\max} \epsilon \zeta}{1 - \gamma} \quad (15)$$

$$= - \frac{V_{\max} \epsilon \zeta}{1 - \gamma} \quad (16)$$

75 The last step follows from the first part in the proof of Lemma 2, $v_{M'}^\pi - v_{M'}^{\Xi(\pi)} \leq \frac{V_{\max} \epsilon \zeta}{1 - \gamma}$. \square

76 B Justification of Assumption 1

77 In this section we prove a claim stated in Section 5 about the upper bound on density functions. We
78 are going to prove Assumption 1 holds under when the transition density is bounded.

79 **Lemma 4.** *Let $p(\cdot|s, a)$ be the probability density function of transition distribution: $\rho(s_0) \leq \sqrt{U} <$
80 ∞ , $p(s_{t+1}|s_t, a_t) \leq \sqrt{U} < \infty$ and $\forall \pi(a_t|s_t, h) \leq \sqrt{U} < \infty$, for all $s_0, s_t, s_{t+1} \in \mathcal{S}$ and $a \in \mathcal{A}$.
81 Then in M' for any non-stationary policy $\pi : \mathcal{S} \times \mathbb{N} \rightarrow \Delta(\mathcal{A}')$ and $h \geq 0$, $\eta_h^\pi(s, a) \leq U$ for any
82 $s \in \mathcal{S}$ and $a \in \mathcal{A}$.*

83 *Proof.* We first prove that $\eta_h^\pi(s) \leq \sqrt{U}$ for any non-stationary policy π . For $h = 0$, $\eta_h^\pi(s) = \rho(s) \leq$
84 \sqrt{U} . For $h \geq 1$ and $s \in \mathcal{S}$:

$$\eta_h^\pi(s) = \int_{s_{-1} \in \mathcal{S}'} \sum_{a \in \mathcal{A}'} \eta_{h-1}^\pi(s_{-1}) \pi(a_{-1}|s_{-1}, h-1) p(s|s_{-1}, a_{-1}) ds_{-1} \quad (17)$$

$$= \int_{s_{-1} \in \mathcal{S}} \sum_{a \in \mathcal{A}} \eta_{h-1}^\pi(s_{-1}) \pi(a_{-1}|s_{-1}, h-1) p(s|s_{-1}, a_{-1}) ds_{-1} \quad (18)$$

$$\leq \mathbb{E}_{\eta_{h-1}^\pi \times \pi(h-1)} [p(s|s_{-1}, a_{-1})] \quad (19)$$

$$\leq \sqrt{U} \quad (20)$$

85 The first step follows from the inductive definition of $\eta_h^\pi(s)$. The second step follows from that s_{abs} is
86 absorbing state and a_{abs} only leads to absorbing state. The third step follows from transition density
87 $p(s|s_{-1}, a_{-1})$ is non-negative. The last step follows from that the transition density $p(s|s_{-1}, a_{-1})$
88 is the same between M and M' for $s, s_{-1} \in \mathcal{S}$, $a_{-1} \in \mathcal{A}$, and $p(s|s_{-1}, a_{-1})$ in M is upper bounded
89 by U . Finally, the joint density function over s and a $\eta_h^\pi(s, a) = \eta_h^\pi(s) \pi(a|s, h)$ is bounded by U ,
90 and we finished the proof. \square

91 For the convenience of notation later we use *admissible distribution* to refer to state-action distribu-
92 tions introduced by non-stationary policy π in M' . This definition is from [1]:

93 **Definition 4** (Admissible distributions). *We say a distribution or its density function $\nu \in \Delta(\mathcal{S}' \times \mathcal{A}')$*
94 *is admissible in MDP M' , if there exists $h \geq 0$, and a (non-stationary) policy $\pi : \mathcal{S}' \times \mathbb{N} \rightarrow \Delta(\mathcal{A}')$,*
95 *such that $\nu(s, a) = \eta_h^\pi(s, a)$.*

96 **C Proofs for Policy Iteration Guarantees**

97 In this section we are going to prove the result of Theorem 1 using the definition of the strong ζ -
 98 constrained policy set . At a high level, the proof is done in two steps. First we prove similar result
 99 to Theorem 1 for any policy in the strong ζ -constrained policy set : an upper bound of $v_{M'}^\pi - v_{M'}^{\pi_t}$,
 100 where π can be any policy in the strong ζ -constrained policy set and π_t is the output of the algorithm
 101 (Theorem 2, formally stated in Appendix C.4). Then we are going to show that for any policy π
 102 in the ζ -constrained policy set after a projection Ξ it is in the strong ζ -constrained policy set and
 103 $v_M^\pi \leq v_M^{\Xi(\pi)} + \frac{V_{\max} \epsilon_\zeta}{1-\gamma}$. Then we can provide the upper bound for $v_M^\pi - v_M^{\pi_t}$ for any π in ζ -constrained
 104 policy set .

105 The proof of Theorem 2 (the Π_{SC}^{all} version of Theorem 1, formally stated in Appendix C.4) goes as
 106 follow. First, we show the fixed point of \mathcal{T}_ζ^π is $Q^{\Xi(\pi)}$ for any policy π , indicating the inner loop of
 107 policy evaluation step is actually evaluating $\pi_t = \Xi(\hat{\pi}_t)$. We prove this result formally in Lemma 6.

108 To bound the gap between π_t and any policy $\tilde{\pi}$ in the ζ -constrained policy set , we use the contraction
 109 property of \mathcal{T}_ζ^π to recursively decompose it into a discounted summation over policy improvement
 110 gap $Q^{\pi_{t+1}} - Q^{\pi_t}$. $\tilde{\pi}$ in the ζ -constrained policy set is needed because the operator \mathcal{T}_ζ^π constrains
 111 the backup on the support set of ζ .

Next, we bound the policy improvement gap in Lemma 12:

$$Q^{\pi_{t+1}} - Q^{\pi_t} \geq -\mathcal{O}(\|\zeta(Q^{\pi_t} - f_{t,K})\|_{1,\nu})$$

112 for some admissible distribution ν related to π_{t+1} . The fact that we only need to measure the error
 113 on the support set of ζ is important. It follows from the fact that both π_{t+1} and π_t only takes action
 114 on the support set of ζ except a_{abs} which gives us a constant value. This allows us to change the
 115 measure from arbitrary distribution ν to data distribution μ , *without needing concentratability*.

116 The rest of proof is to upper bound $\|\zeta(Q^{\pi_t} - f_{t,K})\|_{1,\nu}$ using contraction and concentration inequali-
 117 ties. First, $\|\zeta(Q^{\pi_t} - f_{t,K})\|_{1,\nu}$ is upper bounded by $C\|f_{t,K} - \mathcal{T}_\zeta^\pi f_{t,K}\|_{2,\mu}/(1-\gamma)$ in Lemma 9, using
 118 a standard contraction analysis technique. Notice that here we can change the measure to μ with
 119 cost C to allow us to apply concentration inequality. Then Lemma 8 bounds $\|f_{t,K} - \mathcal{T}_\zeta^\pi f_{t,K}\|_{2,\mu}$
 120 by a function of sample size n and completeness error $\epsilon_{\mathcal{F}}$ using Bernstein's inequality.

121 While writing the proof, we will first introduce the fixed point of \mathcal{T}_ζ^π is $Q^{\Xi(\pi)}$ in section C.1. We
 122 prove the upper bound of the policy evaluation error $\|\zeta(Q^{\pi_t} - f_{t,K})\|_{1,\nu}$, in section C.2, and the
 123 policy improvement step in section C.3. After we proved the main theorem, we will prove when we
 124 can bound the value gap with the optimal value in Corollary 1, as we showed in the main text.

125 **C.1 Fixed point property**

126 In Algorithm 1, the output policy is $\hat{\pi}_{t+1}$. However, we will show that is actually equivalent with
 the following algorithm,

Algorithm 3 Pessimistic Policy Iteration (PPI, repeat Algorithm 1)

Input: $D, \mathcal{F}, \Pi, \hat{\mu}, b$
Output: $\hat{\pi}_T$
 Initialize $\pi_0 \in \Pi$.
for $t = 0$ **to** $T - 1$ **do**
 Initialize $f_{t,0} \in \mathcal{F}$
 for $k = 0$ **to** K **do**
 // Policy Evaluation
 $f_{t,k+1} \leftarrow \arg \min_{f \in \mathcal{F}} \mathcal{L}_D(f, f_{t,k}; \pi_t)$
 end for
 // Policy Improvement
 $\hat{\pi}_{t+1} \leftarrow \arg \max_{\pi \in \Pi} \mathbb{E}_D[\mathbb{E}_\pi [\zeta(s, a) f_{t,K}(s, a)]]$
 $\pi_{t+1} \leftarrow \Xi(\hat{\pi}_{t+1})$
end for

127

128 The output policy is still $\hat{\pi}_{t+1}$, and we know that $v^{\hat{\pi}_{t+1}} \geq v^{\pi_{t+1}}$. So if we can lower bound $v^{\pi_{t+1}}$
129 we immediately have the lower bound on $v^{\hat{\pi}_{t+1}}$. The only difference in algorithm is we change the
130 policy evaluation operator from $\mathcal{T}_\zeta^{\hat{\pi}_t}$ to $\mathcal{T}_\zeta^{\pi_t}$, where π_t is the projection of $\hat{\pi}_t$. The following result
131 shows these two operators are actually the same. For the ease of notation, we refer to Algorithm 3
132 in our analysis.

133 **Lemma 5.** For any policy $\pi : \mathcal{S}' \rightarrow \Delta(\mathcal{A}')$, $\mathcal{T}_\zeta^\pi = \mathcal{T}_\zeta^{\Xi(\pi)}$.

134 *Proof.* We only need to prove for any f , $\mathcal{T}_\zeta^\pi f = \mathcal{T}_\zeta^{\Xi(\pi)} f$. For any $a \in \mathcal{A}$,

$$(\mathcal{T}_\zeta^\pi f)(s, a) = r(s, a) + \gamma \mathbb{E} \left[\sum_{a' \in \mathcal{A}'} \pi(a'|s') \zeta(s', a') f(s', a') \right] \quad (21)$$

$$= r(s, a) + \gamma \mathbb{E} \left[\sum_{a' \in \mathcal{A}'} \pi(a'|s') \zeta^2(s', a') f(s', a') \right] \quad (22)$$

$$= r(s, a) + \gamma \mathbb{E}_{s'} \left[\sum_{a' \in \mathcal{A}'} \Xi(\pi_t)(a'|s') \zeta(s', a') Q^\pi(s', a') \right] \quad (23)$$

$$= (\mathcal{T}_\zeta^{\Xi(\pi)} f)(s, a) \quad (24)$$

135 For $a = a_{\text{abs}}$, $(\mathcal{T}_\zeta^\pi f)(s, a) = 0 = (\mathcal{T}_\zeta^{\Xi(\pi)} f)(s, a)$. \square

136 The next result is a key insight about \mathcal{T}_ζ^π 's behavior in M' that guide our analysis.

137 **Lemma 6.** For any policy $\pi : \mathcal{S}' \rightarrow \Delta(\mathcal{A}')$, the fixed point solution of \mathcal{T}_ζ^π is equal to $Q^{\Xi(\pi)}$ on
138 $\mathcal{S} \times \mathcal{A}$.

139 *Proof.* By definition $Q^{\Xi(\pi)}$ is the fixed point of the standard Bellman evaluation operator on M' :
140 $\mathcal{T}_{M'}^{\Xi(\pi)}$. So for any $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$Q^{\Xi(\pi)}(s, a) \quad (25)$$

$$= (\mathcal{T}_{M'}^{\Xi(\pi)} Q^{\Xi(\pi)})(s, a) \quad (26)$$

$$= r(s, a) + \gamma \mathbb{E}_{s'} \left[\sum_{a' \in \mathcal{A}'} \Xi(\pi)(a'|s') Q^{\Xi(\pi)}(s', a') \right] \quad (27)$$

$$= r(s, a) + \gamma \mathbb{E}_{s'} \left[\Xi(\pi)(a_{\text{abs}}|s') Q^{\Xi(\pi)}(s', a_{\text{abs}}) + \sum_{a' \in \mathcal{A}'} \Xi(\pi)(a'|s') Q^{\Xi(\pi)}(s', a') \right] \quad (28)$$

$$= r(s, a) + \gamma \mathbb{E}_{s'} \left[\sum_{a' \in \mathcal{A}'} \Xi(\pi)(a'|s') Q^{\Xi(\pi)}(s', a') \right] \quad (29)$$

$$= r(s, a) + \gamma \mathbb{E}_{s'} \left[\sum_{a' \in \mathcal{A}'} \pi(a'|s') \zeta(s', a') Q^{\Xi(\pi)}(s', a') \right] \quad (30)$$

$$= (\mathcal{T}_\zeta^\pi Q^{\Xi(\pi)})(s, a) \quad (31)$$

141 So we proved that $Q^{\Xi(\pi)}$ is also the fixed-point solution of \mathcal{T}_ζ^π constrained on $\mathcal{S} \times \mathcal{A}$. \square

142 An obvious consequences of these two lemmas is that the fixed point solution of $\mathcal{T}_\zeta^{\pi_t} = \mathcal{T}_\zeta^{\hat{\pi}_t}$ equals
143 Q^{π_t} on $\mathcal{S} \times \mathcal{A}$.

144 **C.2 Proofs for policy evaluation step**

145 We start with an useful result of the expected loss of the solution from empirical loss minimization,
146 by applying a concentration inequality.

147 **Lemma 7.** *Given $\pi \in \Xi(\Pi)$ and Assumption 3, let $g_f^* = \arg \min_{g \in \mathcal{F}} \|g - \mathcal{T}_\zeta^\pi f\|_{2,\mu}$, then $\|g_f^* -$
148 $\mathcal{T}_\zeta^\pi f\|_{2,\mu}^2 \leq \epsilon_{\mathcal{F}}$. The dataset D is generated i.i.d. from M as follows: $(s, a) \sim \mu$, $r = R(s, a)$,
149 $s' \sim P(s, a)$. Define $\mathcal{L}_\mu(f; f', \pi) = \mathbb{E}_D [\mathcal{L}_D(f; f', \pi)]$. We have that $\forall f \in \mathcal{F}$, with probability at
150 least $1 - \delta$,*

$$\mathcal{L}_\mu(\mathcal{T}_{\zeta,D} f; f, \pi) - \mathcal{L}_\mu(g_f^*; f, \pi) \leq \frac{112V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{3n} + \sqrt{\frac{64V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{n}} \epsilon_{\mathcal{F}}$$

151 where $\mathcal{T}_{\zeta,D}^\pi f = \arg \min_{g \in \mathcal{F}} \mathcal{L}_D(g; f, \pi)$.

152 *Proof.* This proof is similar with the proof of Lemma 16 in [1], and we adapt it to the ζ -constrained
153 Bellman evaluation operator \mathcal{T}_ζ^π . First, there is no difference in \mathcal{L}_D and \mathcal{L}_μ between M and M' ,
154 and the right hand side is also the same constant for M and M' . The distribution of D in M and M'
155 are the same, since μ does not cover s_{abs} and a_{abs} . So we are going to prove the inequality for M ,
156 and thus this bound holds for M' too.

157 For the simplicity of notations, let $V_f^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \zeta(s, a) f(s, a)$. Fix any $f, g \in \mathcal{F}$, and
158 define

$$X(g, f, g_f^*) := (g(s, a) - r - \gamma V_f^\pi(s'))^2 - (g_f^*(s, a) - r - \gamma V_f^\pi(s'))^2. \quad (32)$$

159 Plugging each $(s, a, r, s') \in D$ into $X(g, f, g_f^*)$, we get i.i.d. variables
160 $X_1(g, f, g_f^*), X_2(g, f, g_f^*), \dots, X_n(g, f, g_f^*)$. It is easy to see that

$$\frac{1}{n} \sum_{i=1}^n X_i(g, f, g_f^*) = \mathcal{L}_D(g; f, \pi) - \mathcal{L}_D(g_f^*; f, \pi). \quad (33)$$

161 By the definition of \mathcal{L}_μ , it is also easy to show that

$$\mathcal{L}_\mu(g; f, \pi) = \|g - \mathcal{T}_\zeta^\pi f\|_{2,\mu}^2 + \mathbb{E}_{s,a \sim \mu} \left[\mathbb{V}_{r,s'} \left(r + \gamma \sum_{a' \in \mathcal{A}} \pi(a'|s') \zeta(s', a') f(s', a') \right) \right], \quad (34)$$

162 where $\mathbb{V}_{r,s'}$ is the variance over conditional distribution of r and s' given (s, a) . Notice that the
163 second part does not depends on g . Then

$$\mathcal{L}_\mu(g; f, \pi) - \mathcal{L}_\mu(\mathcal{T}_\zeta^\pi f; f, \pi) = \|g - \mathcal{T}_\zeta^\pi f\|_{2,\mu}^2 \quad (35)$$

164 Then we bound the variance of X :

$$\begin{aligned} \mathbb{V}[X(g, f, g_f^*)] &\leq \mathbb{E}[X(g, f, g_f^*)^2] \\ &= \mathbb{E}_\mu \left[\left((g(s, a) - r - \gamma V_f^\pi(s'))^2 - (g_f^*(s, a) - r - \gamma V_f^\pi(s'))^2 \right)^2 \right] \\ &\hspace{15em} \text{(Definition of } X) \\ &= \mathbb{E}_\mu \left[(g(s, a) - g_f^*(s, a))^2 (g(s, a) + g_f^*(s, a) - 2r - 2\gamma V_f^\pi(s'))^2 \right] \\ &\leq 4V_{\max}^2 \mathbb{E}_\mu \left[(g(s, a) - g_f^*(s, a))^2 \right] \\ &= 4V_{\max}^2 \|g - g_f^*\|_{2,\mu}^2 \\ &\leq 8V_{\max}^2 (\mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F}}). \end{aligned} \quad (36)$$

165 The last step holds because

$$\begin{aligned}
& \|g - g_f^*\|_{2,\mu}^2 \\
& \leq 2 (\|g - \mathcal{T}_\zeta^\pi f\|_{2,\mu}^2 + \|\mathcal{T}_\zeta^\pi f - g_f^*\|_{2,\mu}^2) \quad ((a+b)^2 \leq 2a^2 + 2b^2) \\
& = 2 (\|g - \mathcal{T}_\zeta^\pi f\|_{2,\mu}^2 - \|\mathcal{T}_\zeta^\pi f - g_f^*\|_{2,\mu}^2 + 2\|\mathcal{T}_\zeta^\pi f - g_f^*\|_{2,\mu}^2) \\
& = 2 [\mathcal{L}_\mu(g; f, \pi) - \mathcal{L}_\mu(\mathcal{T}_\zeta^\pi f; f, \pi)] - (\mathcal{L}_\mu(g_f^*; f, \pi) - \mathcal{L}_\mu(\mathcal{T}_\zeta^\pi f; f, \pi)) + 2\|\mathcal{T}_\zeta^\pi f - g_f^*\|_{2,\mu}^2 \\
& \quad \text{(Equation (35))} \\
& = 2 [\mathcal{L}_\mu(g; f, \pi) - \mathcal{L}_\mu(g_f^*; f, \pi) + 2\|\mathcal{T}_\zeta^\pi f - g_f^*\|_{2,\mu}^2] \\
& = 2 (\mathbb{E}[X(g, f, g_f^*)] + 2\|\mathcal{T}_\zeta^\pi f - g_f^*\|_{2,\mu}^2) \\
& \leq 2(\mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F}})
\end{aligned}$$

166 Next, we apply (one-sided) Bernstein's inequality and union bound over all $f \in \mathcal{F}$, $g \in \mathcal{F}$, and
167 $\pi \in \Xi(\Pi)$. With probability at least $1 - \delta$, we have

$$\begin{aligned}
\mathbb{E}[X(g, f, g_f^*)] - \frac{1}{n} \sum_{i=1}^n X_i(f, f, g_f^*) & \leq \sqrt{\frac{2V[X(g, f, g_f^*) \ln \frac{|\mathcal{F}|^2 |\Pi|}{\delta}] + 4V_{\max}^2 \ln \frac{|\mathcal{F}|^2 |\Pi|}{\delta}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}|^2 |\Pi|}{\delta}}{3n} \\
& = \sqrt{\frac{32V_{\max}^2 (\mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F}}) \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{n}} + \frac{8V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{3n}. \quad (37)
\end{aligned}$$

Since $\mathcal{T}_{\zeta,D}^\pi f$ minimizes $\mathcal{L}_D(\cdot; f, \pi)$, it also minimizes $\frac{1}{n} \sum_{i=1}^n X_i(\cdot, f, g_f^*)$. This is because the two objectives only differ by a constant $\mathcal{L}_D(g_f^*; f, \pi)$. Hence,

$$\frac{1}{n} \sum_{i=1}^n X_i(\mathcal{T}_{\zeta,D}^\pi f, f, g_f^*) \leq \frac{1}{n} \sum_{i=1}^n X_i(g_f^*, f, g_f^*) = 0.$$

168 Then,

$$\mathbb{E}[X(\mathcal{T}_{\zeta,D}^\pi f, f, g_f^*)] \leq 0 + \sqrt{\frac{32V_{\max}^2 (\mathbb{E}[X(\mathcal{T}_{\zeta,D}^\pi f, f, g_f^*)] + 2\epsilon_{\mathcal{F}}) \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{n}} + \frac{8V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{3n}.$$

169 Solving for the quadratic formula,

$$\begin{aligned}
\mathbb{E}[X(\mathcal{T}_{\zeta,D}^\pi f, f, g_f^*)] & \leq \sqrt{48 \left(\frac{8V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{3n} \right)^2 + \frac{64V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{n} \epsilon_{\mathcal{F}} + \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{3n}} \\
& \leq \frac{(56 + 32\sqrt{3})V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{3n} + \sqrt{\frac{64V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{n} \epsilon_{\mathcal{F}}} \\
& \quad (\sqrt{a+b} \leq \sqrt{a} + \sqrt{b} \text{ and } \ln \frac{|\mathcal{F}|}{\delta} > 0) \\
& \leq \frac{112V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{3n} + \sqrt{\frac{64V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{n} \epsilon_{\mathcal{F}}}
\end{aligned}$$

170 Noticing that $\mathbb{E}[X(\mathcal{T}_{\zeta,D}^\pi f, f, g_f^*)] = \mathcal{L}_\mu(\mathcal{T}_{\zeta,D}^\pi f; f, \pi) - \mathcal{L}_\mu(g_f^*; f, \pi)$, we complete the proof. \square

171 **Lemma 8** (Policy Evaluation Accuracy). For any $t, k \geq 1$ and π_t , $f_{t,k}$ and $f_{t,k-1}$ from Algorithm
172 I ,

$$\left\| f_{t,k} - \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} \right\|_{2,\mu}^2 \leq \epsilon_1$$

173 where $\epsilon_1 = \frac{208V_{\max}^2 \ln \frac{|\mathcal{F}| |\Pi|}{\delta}}{3n} + 2\epsilon_{\mathcal{F}}$.

Proof.

$$\begin{aligned}
& \left\| f_{t,k} - \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} \right\|_{2,\mu}^2 \\
&= \mathcal{L}_\mu(f_{t,k}; f_{t,k-1}, \pi_t) - \mathcal{L}_\mu(\mathcal{T}_\zeta^{\pi_t} f_{t,k-1}; f_{t,k-1}, \pi_t) \\
&= \left(\mathcal{L}_\mu(f_{t,k}; f_{t,k-1}, \pi_t) - \mathcal{L}_\mu(g_{f_{t,k-1}}^*; f_{t,k-1}, \pi_t) \right) - \left(\mathcal{L}_\mu(\mathcal{T}_\zeta^{\pi_t} f_{t,k-1}; f_{t,k-1}, \pi_t) - \mathcal{L}_\mu(g_{f_{t,k-1}}^*; f_{t,k-1}, \pi_t) \right) \\
&\leq \frac{112V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{3n} + \sqrt{\frac{64V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{n}} \epsilon_{\mathcal{F}} + \left\| g_{f_{t,k-1}}^* - \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} \right\|_{2,\mu} \\
&\hspace{15em} \text{(Equation (35) and Lemma 7)} \\
&\leq \frac{112V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{3n} + \sqrt{\frac{64V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{n}} \epsilon_{\mathcal{F}} + \epsilon_{\mathcal{F}} \quad \text{(Definition of } g_{f_{t,k-1}}^* \text{ and Assumption 3)} \\
&\leq \frac{112V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{3n} + \frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{n} + \epsilon_{\mathcal{F}} + \epsilon_{\mathcal{F}} = \epsilon_1 \quad (\sqrt{2ab} \leq a + b)
\end{aligned}$$

174

□

175 From this lemma to the proof of main theorem, we are going to condition on the fact that the event
176 in Assumption 2 holds. In the proof of the main theorem we will impose the union bound on all
177 failures.

178 **Lemma 9.** For any admissible distribution ν on $\mathcal{S}' \times \mathcal{A}'$, and any π_t from Algorithm 1.

$$\left\| \zeta(s, a) (f_{t,K}(s, a) - Q^{\pi_t}(s, a)) \right\|_{1,\nu} \leq \frac{C(\sqrt{\epsilon_1} + V_{\max}\epsilon_\mu)}{1 - \gamma} + \gamma^K V_{\max} \quad (38)$$

179 where ϵ_1 is defined in Lemma 8.

180 (Although $f_{t,K}$ is only defined on $\mathcal{S} \times \mathcal{A}$, ζ is always zero for any other (s, a) . Thus the all values
181 used in the proof are well-defined. Later, when it is necessary for proof, we define the value of $f_{t,K}$
182 outside of $\mathcal{S} \times \mathcal{A}$ to be zero. In the algorithm, we will never need to query the value of $f_{t,K}$ outside
183 of $\mathcal{S} \times \mathcal{A}$.)

184 *Proof.* For any $k \geq 1$ and any distribution ν on $\mathcal{S}' \times \mathcal{A}'$:

$$\left\| \zeta(f_{t,k} - Q^{\pi_t}) \right\|_{1,\nu} \quad (39)$$

$$\leq \left\| \zeta(f_{t,k} - \mathcal{T}_\zeta^{\pi_t} f_{t,k-1}) \right\|_{1,\nu} + \left\| \zeta(\mathcal{T}_\zeta^{\pi_t} f_{t,k-1} - \mathcal{T}_\zeta^{\pi_t} Q^{\pi_t}) \right\|_{1,\nu} \quad (40)$$

$$\leq \left\| \zeta(f_{t,k} - \mathcal{T}_\zeta^{\pi_t} f_{t,k-1}) \right\|_{1,\nu} + \left\| \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} - \mathcal{T}_\zeta^{\pi_t} Q^{\pi_t} \right\|_{1,\nu} \quad (41)$$

$$\leq C \left\| f_{t,k} - \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} \right\|_{1,\hat{\mu}} + \left\| \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} - \mathcal{T}_\zeta^{\pi_t} Q^{\pi_t} \right\|_{1,\nu} \quad (42)$$

$$\leq C \left(\left\| f_{t,k} - \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} \right\|_{1,\mu} + V_{\max}\epsilon_\mu \right) + \left\| \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} - \mathcal{T}_\zeta^{\pi_t} Q^{\pi_t} \right\|_{1,\nu} \quad (43)$$

$$\leq C \left(\left\| f_{t,k} - \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} \right\|_{2,\mu} + V_{\max}\epsilon_\mu \right) + \left\| \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} - \mathcal{T}_\zeta^{\pi_t} Q^{\pi_t} \right\|_{1,\nu} \quad \text{(Jensen's inequality)}$$

$$\leq C(\sqrt{\epsilon_1} + V_{\max}\epsilon_\mu) + \left\| \mathcal{T}_\zeta^{\pi_t} f_{t,k-1} - \mathcal{T}_\zeta^{\pi_t} Q^{\pi_t} \right\|_{1,\nu} \quad \text{(Lemma 8)}$$

$$= C(\sqrt{\epsilon_1} + V_{\max}\epsilon_\mu) + \mathbb{E}_\nu \left| \gamma \mathbb{E}_{P(\nu)} \sum_{a' \in \mathcal{A}} \pi_t(a' | s') \zeta(s', a') (f_{t,k-1}(s', a') - Q^{\pi_t}(s', a')) \right| \quad (44)$$

$$= C(\sqrt{\epsilon_1} + V_{\max}\epsilon_\mu) + \mathbb{E}_\nu \left[\gamma \mathbb{E}_{P(\nu) \times \pi_t} |\zeta(s', a') (f_{t,k-1}(s', a') - Q^{\pi_t}(s', a'))| \right] \quad (45)$$

$$\leq C(\sqrt{\epsilon_1} + V_{\max}\epsilon_\mu) + \gamma \mathbb{E}_{P(\nu) \times \pi_t} |\zeta(s', a') (f_{t,k-1}(s', a') - Q^{\pi_t}(s', a'))| \quad (46)$$

$$\leq C(\sqrt{\epsilon_1} + V_{\max}\epsilon_\mu) + \gamma \left\| \zeta(f_{t,k-1} - Q^{\pi_t}) \right\|_{1,P(\nu) \times \pi} \quad (47)$$

185 Equation (42) holds since for all (s, a) s.t. $\zeta(s, a) > 0$, $\nu(s, a) \leq U \leq \frac{U}{b} \widehat{\mu}(s, a) = C \widehat{\mu}(s, a)$.
 186 Equation (43) holds since the total variation distance between μ and $\widehat{\mu}$ is bounded by ϵ_μ and the
 187 Bellman error is bounded in $[-V_{\max}, V_{\max}]$. Equation (44) follows from $\pi_t \in \Pi_{SC}^{all}$. So if $\zeta(s, a) =$
 188 0 , $\pi(a|s) = 0$ for all $a \in \mathcal{A}$. Equation (45) holds since $\zeta(\cdot, a_{\text{abs}}) = 0$. The next equation follows
 189 from that $\zeta = \zeta^2$.

190 Note that this holds for any admissible distribution ν on $\mathcal{S}' \times \mathcal{A}'$ and and k , as well as ϵ_1 does not
 191 depends on k . Repeating this for k from K to 1 we will have that

$$\|\zeta(s, a)(f_{t,K}(s, a) - Q^{\pi_t}(s, a))\|_{1,\nu} \leq \frac{1 - \gamma^K}{1 - \gamma} C (\sqrt{\epsilon_1} + V_{\max} \epsilon_\mu) + \gamma^K V_{\max} \quad (48)$$

$$< \frac{C (\sqrt{\epsilon_1} + V_{\max} \epsilon_\mu)}{1 - \gamma} + \gamma^K V_{\max} \quad (49)$$

192

□

193 C.3 Proofs for policy improvement step

Lemma 10 (Concentration of Policy Improvement Loss). *For any $f \in \mathcal{F}$, with probability at least $1 - \delta$,*

$$\left\| \mathbb{E}_{\widehat{\pi}_f} [\zeta(s, a)f(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a)f(s, a) \right\|_{1,\mu} \leq \epsilon_\Pi + 2V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}}$$

194 where $\widehat{\pi}_f = \arg \max_{\pi \in \Pi} \mathbb{E}_D [\mathbb{E}_\pi [\zeta(s, a)f(s, a)]]$.

195 *Proof.* Fixed f , define $X(s; \pi) = \max_{a \in \mathcal{A}} \zeta(s, a)f(s, a) - \mathbb{E}_\pi [\zeta(s, a)f(s, a)]$. Notice that by
 196 definition $X(s; \pi)$ is always non-negative, and $\widehat{\pi}_f = \arg \max_{\pi \in \Pi} \mathbb{E}_D [\mathbb{E}_\pi [\zeta(s, a)f(s, a)]] =$
 197 $\arg \min_{\pi \in \Pi} \mathbb{E}_D [X(s; \pi)]$.

Only in this proof, let π_f be:

$$\arg \min_{\pi \in \Pi} \mathbb{E}_\mu [X(s; \pi)] = \arg \min_{\pi \in \Pi} \left\| \mathbb{E}_\pi [\zeta(s, a)f(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a)f(s, a) \right\|_{1,\mu}.$$

198 $X(s; \pi) \in [0, V_{\max}]$. By Hoeffding's inequality and union bound over all $\pi \in \Pi$, $f \in \mathcal{F}$, with
 199 probability at least $1 - \delta$ for any f and $\pi \neq \pi_f$,

$$\mathbb{E}_\mu [X(s; \pi)] - \mathbb{E}_D [X(s; \pi)] \leq V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \quad (50)$$

200 for $\pi = \pi_f$

$$\mathbb{E}_D [X(s; \pi)] - \mathbb{E}_\mu [X(s; \pi)] \leq V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \quad (51)$$

201 If $\widehat{\pi}_f = \pi_f$, then $\mathbb{E}_\mu [X(s; \widehat{\pi}_f)] \leq \epsilon_\Pi$. Otherwise,

$$\mathbb{E}_\mu [X(s; \widehat{\pi}_f)] \quad (52)$$

$$\leq \mathbb{E}_D [X(s; \widehat{\pi}_f)] + V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \quad (53)$$

$$\leq \mathbb{E}_D [X(s; \pi_f)] + V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \quad (54)$$

$$\leq \mathbb{E}_\mu [X(s; \pi_f)] + 2V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \quad (55)$$

$$= \min_{\pi \in \Pi} \left\| \mathbb{E}_{\widehat{\pi}} [\zeta(s, a)f(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a)f(s, a) \right\|_{1,\mu} + 2V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \quad (56)$$

$$= \epsilon_\Pi + 2V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \quad (57)$$

202

□

203 For the following proof until the main theorem, we are going to condition on the fact that the high
 204 probability bound in the lemma above holds, and impose an union bound in the proof of main
 205 theorem.

206 **Lemma 11.** For any admissible distribution ν on S' , any policy $\pi : S' \rightarrow \Delta(\mathcal{A}')$,

$$\begin{aligned} & \mathbb{E}_\nu [\mathbb{E}_{\pi_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \mathbb{E}_\pi [\zeta(s, a) f_{t,K}(s, a)]] \geq \\ & -C \left(\epsilon_\Pi + V_{\max} \epsilon_\mu + 2V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \right) \end{aligned}$$

207 *Proof.* Recall that $\pi_{t+1} = \Xi(\hat{\pi}_{t+1})$. So $\pi_{t+1}(a|s) = \hat{\pi}_{t+1}(a|s)$ for all a such that $\zeta(s, a) = 1$. Then

$$\begin{aligned} & \mathbb{E}_{\pi_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] = \mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] \\ & \mathbb{E}_\nu [\mathbb{E}_{\pi_{t+1}} [\zeta(s, a) f_{t,K}(s, a)]] = \mathbb{E}_\nu [\mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)]] \end{aligned}$$

$$\mathbb{E}_\nu [\mathbb{E}_{\pi_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \mathbb{E}_\pi [\zeta(s, a) f_{t,K}(s, a)]] \quad (58)$$

$$= \mathbb{E}_\nu [\mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \mathbb{E}_\pi [\zeta(s, a) f_{t,K}(s, a)]] \quad (59)$$

$$= \mathbb{E}_\nu \left[\mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) + \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) - \mathbb{E}_\pi [\zeta(s, a) f_{t,K}(s, a)] \right] \quad (60)$$

$$\geq \mathbb{E}_\nu \left[\mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) \right] \quad (61)$$

$$\geq -\mathbb{E}_\nu \left| \mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) \right| \quad (62)$$

$$= -\left\| \mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) \right\|_{1, \nu} \quad (63)$$

$$\geq -C \left\| \mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) \right\|_{1, \hat{\mu}} \quad (64)$$

208 The last step follows from that $\zeta(s, a) = 1 \Rightarrow \hat{\mu}(s, a) \geq b \Rightarrow \hat{\mu}(s) \geq b \Rightarrow -\nu(s) \geq -U \geq$
 209 $-C\hat{\mu}(s)$, and for all other (s, a) the term inside of norm is zero. Since the total variation distance
 210 between $\hat{\mu}$ and μ is bounded by ϵ_μ

$$\left\| \mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) \right\|_{1, \hat{\mu}} \quad (65)$$

$$\leq \left\| \mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) \right\|_{1, \mu} + V_{\max} \epsilon_\mu \quad (66)$$

211 By Lemma 10:

$$\left\| \mathbb{E}_{\hat{\pi}_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \max_{a \in \mathcal{A}} \zeta(s, a) f_{t,K}(s, a) \right\|_{1, \mu} \leq \epsilon_\Pi + 2V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \quad (67)$$

212 Then we finished the proof by plug this into the last equation. \square

213 **Lemma 12.** For any $(s, a) \in S' \times \mathcal{A}'$, and any π_t, π_{t+1} in Algorithm 1,

$$Q^{\pi_{t+1}}(s, a) - Q^{\pi_t}(s, a) \geq -\frac{2C\sqrt{\epsilon_1} + 3V_{\max}C\epsilon_\mu}{(1-\gamma)^2} - \frac{\epsilon_2 + 2\gamma^K V_{\max}}{1-\gamma} \quad (68)$$

214 where ϵ_1 is defined in Lemma 8, $\epsilon_2 = C \left(\epsilon_\Pi + 2V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \right)$.

215 *Proof.* For any s' , only in this proof, let $\eta_h^{\pi_{t+1}}$ be the state distribution on the h th step from initial
 216 state s' following π_{t+1} . By applying performance difference lemma [3],

$$V^{\pi_{t+1}}(s') - V^{\pi_t}(s') \tag{69}$$

$$= \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{z \sim \eta_h^{\pi_{t+1}}} \left[\sum_{a \in \mathcal{A}'} (\pi_{t+1}(a|z) Q^{\pi_t}(z, a) - \pi_t(a|z) Q^{\pi_t}(z, a)) \right] \tag{70}$$

$$= \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{z \sim \eta_h^{\pi_{t+1}}} \left[\sum_{a \in \mathcal{A}'} (1 - \zeta(z, a)) (\pi_{t+1}(a|z) Q^{\pi_t}(z, a) - \pi_t(a|z) Q^{\pi_t}(z, a)) \right] \tag{71}$$

$$+ \sum_{a \in \mathcal{A}'} \zeta(z, a) (\pi_{t+1}(a|z) Q^{\pi_t}(z, a) - \pi_t(a|z) Q^{\pi_t}(z, a)) \tag{72}$$

217 Because $\pi_t, \pi_{t+1} \in \Pi_{SC}^{all}$, $\zeta(z, a) = 0$ means either $\pi_t(a|z) = \pi_{t+1}(a|z) = 0$ or $a = a_{\text{abs}}$. So the
 218 first term is zero. Then:

$$V^{\pi_{t+1}}(s') - V^{\pi_t}(s') \tag{73}$$

$$= \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{z \sim \eta_h^{\pi_{t+1}}} \left[\sum_{a \in \mathcal{A}'} \zeta(z, a) (\pi_{t+1}(a|z) Q^{\pi_t}(z, a) - \pi_t(a|z) Q^{\pi_t}(z, a)) \right] \tag{74}$$

$$= \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{z \sim \eta_h^{\pi_{t+1}}} \left[\sum_{a \in \mathcal{A}} \zeta(z, a) (\pi_{t+1}(a|z) Q^{\pi_t}(z, a) - \pi_t(a|z) Q^{\pi_t}(z, a)) \right] \tag{75}$$

$$= \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{z \sim \eta_h^{\pi_{t+1}}} \left[\sum_{a \in \mathcal{A}} \zeta(z, a) (\pi_{t+1}(a|z) Q^{\pi_t}(z, a) - \pi_{t+1}(a|z) f_{t,K}(z, a)) \right] \tag{76}$$

$$+ \sum_{a \in \mathcal{A}} \zeta(z, a) (\pi_{t+1}(a|z) f_{t,K}(z, a) - \pi_t(a|z) f_{t,K}(z, a)) \tag{77}$$

$$+ \sum_{a \in \mathcal{A}} \zeta(z, a) (\pi_t(a|z) f_{t,K}(z, a) - \pi_t(a|z) Q^{\pi_t}(z, a)) \tag{78}$$

219 Equation 75 follows from $Q^\pi(s, a_{\text{abs}}) = 0$ for any π and s . By Lemma 11, for any h ,

$$\mathbb{E}_{z \sim \eta_h^{\pi_{t+1}}} \left[\sum_{a \in \mathcal{A}} \zeta(z, a) (\pi_{t+1}(a|z) f_{t,K}(z, a) - \pi_t(a|z) f_{t,K}(z, a)) \right] \tag{79}$$

$$= \mathbb{E}_{z \sim \eta_h^{\pi_{t+1}}} [\mathbb{E}_{\pi_{t+1}} [\zeta(s, a) f_{t,K}(s, a)] - \mathbb{E}_{\pi_t} [\zeta(s, a) f_{t,K}(s, a)]] \geq -\epsilon_2 - CV_{\max} \epsilon_\mu \tag{80}$$

220 Then

$$V^{\pi_{t+1}}(s') - V^{\pi_t}(s') \quad (81)$$

$$\geq \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{z \sim \eta_h^{\pi_{t+1}}} \left[\sum_{a \in \mathcal{A}} \zeta(z, a) (\pi_{t+1}(a|z) Q^{\pi_t}(z, a) - \pi_{t+1}(a|z) f_{t,K}(z, a)) \right] \quad (82)$$

$$+ \sum_{a \in \mathcal{A}} \zeta(z, a) (\pi_t(a|z) f_{t,K}(z, a) - \pi_t(a|z) Q^{\pi_t}(z, a)) \Big] - \frac{\epsilon_2 + CV_{\max} \epsilon_\mu}{1 - \gamma} \quad (83)$$

$$\geq - \sum_{h=1}^{\infty} \gamma^{h-1} \left(\|\zeta(z, a)(Q^{\pi_t}(z, a) - f_{t,K}(z, a))\|_{1, \eta_h^{\pi_{t+1}}} \right. \quad (84)$$

$$\left. + \|\zeta(z, a)(Q^{\pi_t}(z, a) - f_{t,K}(z, a))\|_{1, \eta_h^{\pi_{t+1}} \times \pi_t} \right) - \frac{\epsilon_2 + CV_{\max} \epsilon_\mu}{1 - \gamma} \quad (85)$$

$$\geq - \sum_{h=1}^{\infty} \gamma^{h-1} \left(\|\zeta(z, a)(Q^{\pi_t}(z, a) - f_{t,K}(z, a))\|_{2, \eta_h^{\pi_{t+1}}} \right. \quad (86)$$

$$\left. + \|\zeta(z, a)(Q^{\pi_t}(z, a) - f_{t,K}(z, a))\|_{2, \eta_h^{\pi_{t+1}} \times \pi_t} \right) - \frac{\epsilon_2 + CV_{\max} \epsilon_\mu}{1 - \gamma} \quad (87)$$

$$\geq \frac{-2C(\sqrt{\epsilon_1} + V_{\max} \epsilon_\mu)}{(1 - \gamma)^2} - \frac{2\gamma^K V_{\max}}{1 - \gamma} - \frac{\epsilon_2 + CV_{\max} \epsilon_\mu}{1 - \gamma} \quad (\text{Lemma 9})$$

221 Equation 87 follows from Jensen's inequality. Since this holds for any s' , we proved that for any
222 (s, a) ,

$$[Q^{\pi_{t+1}}(s, a) - Q^{\pi_t}(s, a)] \quad (88)$$

$$= \gamma \mathbb{E}_{s'} [V^{\pi_{t+1}}(s') - V^{\pi_t}(s')] \quad (89)$$

$$\geq \frac{-2C(\sqrt{\epsilon_1} + V_{\max} \epsilon_\mu)}{(1 - \gamma)^2} - \frac{2\gamma^K V_{\max}}{1 - \gamma} - \frac{\epsilon_2 + CV_{\max} \epsilon_\mu}{1 - \gamma} \quad (90)$$

$$\geq - \frac{2C\sqrt{\epsilon_1} + 3CV_{\max} \epsilon_\mu}{(1 - \gamma)^2} - \frac{2\gamma^K V_{\max}}{1 - \gamma} - \frac{\epsilon_2}{1 - \gamma} \quad (91)$$

223

□

224 C.4 Proof of main theorems

225 **Theorem 2.** Given an MDP $M = \langle \mathcal{S}, \mathcal{A}, R, P, \gamma, p \rangle$, a dataset $D = \{(s, a, r, s')\}$ with n samples
226 that is draw i.i.d. from $\mu \times R \times P$, and a finite Q -function classes \mathcal{F} and a finite policy class Π
227 satisfying Assumption 3 and 4, $\pi_t = \Xi(\hat{\pi}_t)$ from Algorithm 1 satisfies that with probability at least
228 $1 - 3\delta$,

$$v^{\tilde{\pi}} - v^{\pi_t} \leq \frac{4C}{(1 - \gamma)^3} \left(\sqrt{\frac{419V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{3n}} + 2\sqrt{\epsilon_{\mathcal{F}}} \right) + \frac{6CV_{\max} \epsilon_\mu}{(1 - \gamma)^3} + \frac{2C\epsilon_\Pi + 3\gamma^{K-1}V_{\max}}{(1 - \gamma)^2}$$

229 for any policy $\tilde{\pi} \in \Pi_{SC}^{all}$.

230 *Proof.* For simplicity of the notation, let $\epsilon_1 = \frac{208V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{3n} + 2\epsilon_{\mathcal{F}}$, $\epsilon_2 =$
231 $C \left(\epsilon_\Pi + 2V_{\max} \sqrt{\frac{\ln(|\mathcal{F}||\Pi|/\delta)}{2n}} \right)$ and $\epsilon_3 = \frac{2C\sqrt{\epsilon_1} + 3V_{\max} C \epsilon_\mu}{(1 - \gamma)^2} + \frac{\epsilon_2 + 2\gamma^K V_{\max}}{1 - \gamma}$. We start by proving

232 a stronger result. For any $\tilde{\pi} \in \Pi_{SC}^{all}$, we will upper bound $\mathbb{E}_\nu [V^{\tilde{\pi}} - V^{\pi_t}]$ for any admissible

233 distribution ν over S' which will naturally be an upper bound for $v^{\tilde{\pi}} - v^{\pi_t}$

$$\begin{aligned}
& \mathbb{E}_\nu[V^{\tilde{\pi}} - V^{\pi_{t+1}}] \\
&= \mathbb{E}_\nu \left[V^{\tilde{\pi}}(s) - \sum_{a \in \mathcal{A}'} \pi_{t+1}(a|s) Q^{\pi_t}(s, a) + \sum_{a \in \mathcal{A}'} \pi_{t+1}(a|s) Q^{\pi_t}(s, a) - V^{\pi_{t+1}}(s) \right] \\
&= \mathbb{E}_\nu \left[V^{\tilde{\pi}}(s) - \sum_{a \in \mathcal{A}'} \pi_{t+1}(a|s) Q^{\pi_t}(s, a) + \sum_{a \in \mathcal{A}'} \pi_{t+1}(a|s) (Q^{\pi_t}(s, a) - Q^{\pi_{t+1}}(s, a)) \right] \\
&\leq \mathbb{E}_\nu \sum_{a \in \mathcal{A}'} \left[\tilde{\pi}(a|s) Q^{\tilde{\pi}}(s, a) - \pi_{t+1}(a|s) Q^{\pi_t}(s, a) \right] + \epsilon_3 \quad (\text{Lemma 12}) \\
&= \mathbb{E}_\nu \sum_{a \in \mathcal{A}'} \zeta(s, a) [\tilde{\pi}(a|s) Q^{\tilde{\pi}}(s, a) - \pi_{t+1}(a|s) Q^{\pi_t}(s, a)] + \epsilon_3 \\
&= \mathbb{E}_\nu \left[\mathbb{E}_{\tilde{\pi}} \left[\zeta(s, a) Q^{\tilde{\pi}}(s, a) \right] - \mathbb{E}_{\pi_{t+1}} \left[\zeta(s, a) f_t(s, a) \right] \right. \\
&\quad \left. + \mathbb{E}_{\pi_{t+1}} \left[\zeta(s, a) f_t(s, a) \right] - \mathbb{E}_{\pi_{t+1}} \left[\zeta(s, a) Q^{\pi_{t+1}}(s, a) \right] \right] + \epsilon_3 \\
&\leq \mathbb{E}_\nu \left[\mathbb{E}_{\tilde{\pi}} \left[\zeta(s, a) Q^{\tilde{\pi}}(s, a) \right] - \mathbb{E}_{\pi_{t+1}} \left[\zeta(s, a) f_t(s, a) \right] \right] \\
&\quad + \|\zeta(z, a)(Q^{\pi_t}(z, a) - f_t(z, a))\|_{1, \nu \times \pi_{t+1}} + \epsilon_3 \\
&\leq \mathbb{E}_\nu \left[\mathbb{E}_{\tilde{\pi}} \left[\zeta(s, a) Q^{\tilde{\pi}}(s, a) \right] - \mathbb{E}_{\pi_{t+1}} \left[\zeta(s, a) f_t(s, a) \right] \right] + \frac{C\sqrt{\epsilon_1} + CV_{\max}\epsilon_\mu}{1-\gamma} + \gamma^K V_{\max} + \epsilon_3 \\
&\quad (\text{Lemma 9}) \\
&\leq \mathbb{E}_\nu \left[\mathbb{E}_{\tilde{\pi}} \left[\zeta(s, a) Q^{\tilde{\pi}}(s, a) \right] - \mathbb{E}_{\tilde{\pi}} \left[\zeta(s, a) f_t(s, a) \right] \right] + \epsilon_2 + CV_{\max}\epsilon_\mu + \frac{C\sqrt{\epsilon_1} + CV_{\max}\epsilon_\mu}{1-\gamma} + \gamma^K V_{\max} + \epsilon_3 \\
&\quad (\text{Lemma 11}) \\
&\leq \mathbb{E}_\nu \left[\mathbb{E}_{\tilde{\pi}} \left[\zeta(s, a) Q^{\tilde{\pi}}(s, a) \right] - \mathbb{E}_{\tilde{\pi}} \left[\zeta(s, a) Q^{\pi_t}(s, a) \right] \right] + \epsilon_2 + \frac{2C\sqrt{\epsilon_1} + 3CV_{\max}\epsilon_\mu}{1-\gamma} + 2\gamma^K V_{\max} + \epsilon_3 \\
&\quad (\text{Lemma 9}) \\
&= \mathbb{E}_{\nu \times \tilde{\pi}} \left[\zeta(s, a) Q^{\tilde{\pi}}(s, a) - \zeta(s, a) Q^{\pi_t}(s, a) \right] + \epsilon_2 + \frac{2C\sqrt{\epsilon_1} + 3CV_{\max}\epsilon_\mu}{1-\gamma} + 2\gamma^K V_{\max} + \epsilon_3 \\
&= \mathbb{E}_{\nu \times \tilde{\pi}} \left[Q^{\tilde{\pi}}(s, a) - Q^{\pi_t}(s, a) \right] + \epsilon_2 + \frac{2C\sqrt{\epsilon_1} + 3CV_{\max}\epsilon_\mu}{1-\gamma} + 2\gamma^K V_{\max} + \epsilon_3 \quad (\pi_t \in \Pi_{SC}^{all}) \\
&\leq \gamma \mathbb{E}_{P(\nu \times \tilde{\pi})} [V^{\tilde{\pi}} - V^{\pi_t}] + \epsilon_2 + \frac{2C\sqrt{\epsilon_1} + 3CV_{\max}\epsilon_\mu}{1-\gamma} + 2\gamma^K V_{\max} + \epsilon_3
\end{aligned}$$

234 The second to last step follows from $\pi_t \in \Pi_{SC}^{all}$: for all s, a such that $\tilde{\pi}(a|s) > 0$, either $\zeta(s, a) = 1$,
235 or $a = a_{\text{abs}}$. The later two indicate that $Q^{\pi_t}(s, a) = Q^{\tilde{\pi}}(s, a) = 0$. So for all s, a such that
236 $\tilde{\pi}(a|s) > 0$, $Q^{\tilde{\pi}}(s, a) = \zeta(s, a) Q^{\tilde{\pi}}(s, a)$ and $Q^{\pi_t}(s, a) = \zeta(s, a) Q^{\pi_t}(s, a)$. Now we proved

$$\mathbb{E}_\nu[V^{\tilde{\pi}} - V^{\pi_{t+1}}] \leq \gamma \mathbb{E}_{P(\nu \times \tilde{\pi})} [V^{\tilde{\pi}} - V^{\pi_t}] + \epsilon_2 + \epsilon_3 + \frac{2C\sqrt{\epsilon_1} + 3CV_{\max}\epsilon_\mu}{1-\gamma} + 2\gamma^K V_{\max} \quad (92)$$

237 holds for any distribution ν . The error terms do not depend on t and this holds for any t . We can
 238 repeatedly apply this for all $0 < t' \leq t$. Assuming $t \geq K$ this will give us :

$$\begin{aligned}
 & \mathbb{E}_\nu[V^{\tilde{\pi}} - V^{\pi_{t+1}}] \\
 & \leq \frac{1-\gamma^t}{1-\gamma} \left(\epsilon_2 + \epsilon_3 + \frac{2C\sqrt{\epsilon_1} + 3CV_{\max}\epsilon_\mu}{1-\gamma} + 2\gamma^K V_{\max} \right) + \gamma^t V_{\max} \\
 & \leq \frac{\epsilon_2}{1-\gamma} + \frac{\epsilon_3}{1-\gamma} + \frac{2C\sqrt{\epsilon_1}}{(1-\gamma)^2} + \frac{3CV_{\max}\epsilon_\mu}{(1-\gamma)^2} + \frac{3\gamma^K V_{\max}}{1-\gamma} \\
 & \leq \frac{2\epsilon_2}{(1-\gamma)^2} + \frac{4C\sqrt{\epsilon_1}}{(1-\gamma)^3} + \frac{6CV_{\max}\epsilon_\mu}{(1-\gamma)^3} + \frac{3\gamma^{K-1}V_{\max}}{(1-\gamma)^2} \\
 & \leq \frac{2C\epsilon_\Pi}{(1-\gamma)^2} + \frac{4C}{(1-\gamma)^2} \sqrt{\frac{V_{\max}^2 \ln(|\mathcal{F}||\Pi|/\delta)}{2n}} + \frac{4C\sqrt{\epsilon_1}}{(1-\gamma)^3} + \frac{6CV_{\max}\epsilon_\mu}{(1-\gamma)^3} + \frac{3\gamma^{K-1}V_{\max}}{(1-\gamma)^2} \\
 & \leq \frac{2C\epsilon_\Pi}{(1-\gamma)^2} + \frac{4C}{(1-\gamma)^3} \left(\sqrt{\frac{V_{\max}^2 \ln(|\mathcal{F}||\Pi|/\delta)}{2n}} + \sqrt{\frac{208V_{\max}^2 \ln(|\mathcal{F}||\Pi|/\delta)}{3n}} + 2\epsilon_{\mathcal{F}} \right) \\
 & \quad + \frac{6CV_{\max}\epsilon_\mu}{(1-\gamma)^3} + \frac{3\gamma^{K-1}V_{\max}}{(1-\gamma)^2} \\
 & \leq \frac{2C\epsilon_\Pi}{(1-\gamma)^2} + \frac{4C}{(1-\gamma)^3} \left(\sqrt{\frac{V_{\max}^2 \ln(|\mathcal{F}||\Pi|/\delta)}{2n}} + \sqrt{\frac{208V_{\max}^2 \ln(|\mathcal{F}||\Pi|/\delta)}{3n}} + \sqrt{2\epsilon_{\mathcal{F}}} \right) \\
 & \quad + \frac{6CV_{\max}\epsilon_\mu}{(1-\gamma)^3} + \frac{3\gamma^{K-1}V_{\max}}{(1-\gamma)^2} \\
 & \leq \frac{2C\epsilon_\Pi}{(1-\gamma)^2} + \frac{4C}{(1-\gamma)^3} \left(\sqrt{\frac{419V_{\max}^2 \ln(|\mathcal{F}||\Pi|/\delta)}{3n}} + \sqrt{2\epsilon_{\mathcal{F}}} \right) + \frac{6CV_{\max}\epsilon_\mu}{(1-\gamma)^3} + \frac{3\gamma^{K-1}V_{\max}}{(1-\gamma)^2}
 \end{aligned}$$

239 The last step follows from that $a + b \leq \sqrt{2(a^2 + b^2)}$. The error bound is finished by simplifying the
 240 expression. The failure probability 3δ is from the union bound of probability δ on which Assumption
 241 2 fails, probability δ on which Lemma 7 fails, and the probability δ on which Lemma 10 fails. \square

242 Now we are going to use the fact that there is an almost no-value-loss projection from the ζ -
 243 constrained policy set to the strong ζ -constrained policy set in order to prove an error bound w.r.t
 244 any $\tilde{\pi} \in \Pi_C^{all}$.

245 **Theorem 1.** Given an MDP $M = \langle \mathcal{S}, \mathcal{A}, R, P, \gamma, p \rangle$, a dataset $D = \{(s, a, r, s')\}$ with n samples
 246 that is draw i.i.d. from $\mu \times R \times P$, and a finite Q -function classes \mathcal{F} and a finite policy class Π
 247 satisfying Assumption 3 and 4, $\hat{\pi}_t$ from Algorithm 1 satisfies that with probability at least $1 - 3\delta$,

$$v_M^{\tilde{\pi}} - v_M^{\hat{\pi}_t} \leq \frac{4C}{(1-\gamma)^3} \left(\sqrt{\frac{419V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{3n}} + 2\sqrt{\epsilon_{\mathcal{F}}} \right) + \frac{6CV_{\max}\epsilon_\mu}{(1-\gamma)^3} + \frac{2C\epsilon_\Pi + 3\gamma^{K-1}V_{\max}}{(1-\gamma)^2} + \frac{V_{\max}\epsilon_\zeta}{1-\gamma}$$

248 for any policy $\tilde{\pi} \in \Pi_C^{all}$ and only take action over \mathcal{A} .

249 *Proof.* For any policy $\tilde{\pi}$ that only take action over \mathcal{A} , Lemma 3 tells that $v_M^{\tilde{\pi}} \leq v_{M'}^{\Xi(\tilde{\pi})} + \frac{V_{\max}\epsilon_\zeta}{1-\gamma}$.
 250 Since $\pi_t = \Xi(\hat{\pi}_t)$ and $\hat{\pi}_t$ only takes action in \mathcal{A} , by Lemma 1 and Lemma 2 $v_M^{\hat{\pi}_t} = v_{M'}^{\hat{\pi}_t} \geq v_M^{\pi_t}$.
 251 Then $v_M^{\tilde{\pi}} - v_M^{\hat{\pi}_t} \leq v_{M'}^{\Xi(\tilde{\pi})} - v_{M'}^{\hat{\pi}_t} + \frac{V_{\max}\epsilon_\zeta}{1-\gamma}$ and Theorem 2 completes the proof. \square

252 When there exist an optimal policy that is supported well by μ . We can derive the following result
 253 about value gap between learned policy and optimal policy immediately from the main theorem
 254 about approximate policy iteration.

255 **Corollary 2.** *If there exists an π^* on M such that $\Pr(\mu(s, a) \leq 2b|\pi^*) \leq \epsilon$. then under the*
 256 *assumptions of Theorem 1, $\hat{\pi}_t$ from Algorithm 1 satisfies that with probability at least $1 - 3\delta$,*

$$v_M^{\pi^*} - v_M^{\hat{\pi}_t} \leq \frac{4C}{(1-\gamma)^3} \left(\sqrt{\frac{419V_{\max}^2 \ln \frac{|\mathcal{F}||\Pi|}{\delta}}{3n}} + 2\sqrt{\epsilon_{\mathcal{F}}} \right) + \frac{6CV_{\max}\epsilon_{\mu}}{(1-\gamma)^3} \\ + \frac{2C\epsilon_{\Pi} + 3\gamma^{K-1}V_{\max}}{(1-\gamma)^2} + \frac{V_{\max}(\epsilon + C\epsilon_{\mu})}{1-\gamma}$$

257 *Proof.* Given the condition of π^* ,

$$\Pr(\hat{\mu}(s, a) \leq b|\pi^*) \leq \Pr(\mu(s, a) \leq 2b|\pi^*) + \Pr(|\mu(s, a) - \hat{\mu}(s, a)| \geq b|\pi^*) \quad (93)$$

$$\leq \epsilon + \Pr(|\mu(s, a) - \hat{\mu}(s, a)| \geq b|\pi^*) \quad (94)$$

$$\leq \epsilon + \frac{\mathbb{E}_{\eta^{\pi^*}} [|\mu(s, a) - \hat{\mu}(s, a)|]}{b} \quad (95)$$

$$\leq \epsilon + \frac{Ud_{\text{TV}}(\mu(s, a), \hat{\mu}(s, a))}{b} \quad (96)$$

$$\leq \epsilon + C\epsilon_{\mu} \quad (97)$$

258 Then $\pi^* \in \Pi_C^{\text{all}}$ with $\epsilon_{\zeta} = \epsilon + C\epsilon_{\mu}$, and applying Theorem 1 finished the proof. \square

259 C.5 Safe Policy Improvement Result

260 In many scenarios we aim to have a policy improvement that is guaranteed to be no worse than
 261 the data collection policy, which is called safe policy improvement. By abusing the notation a bit,
 262 let $\mu(a|s)$ be a policy that generate the data set. For our algorithm, the safe policy improvement
 263 will hold if $\mu \in \Pi_C^{\text{all}}$. To show $\mu \in \Pi_C^{\text{all}}$, we only need that $\Pr(\mu(s, a) \leq b|\mu) \leq \epsilon_{\zeta}$. When
 264 the state-action space is finite, there must exist an minimum value for all non-zero $\mu(s, a)$'s. Let
 265 $\mu_{\min} = \min_{s, a, s.t. \mu(s, a) > 0} \mu(s, a)$. Then we have that, if $b \leq \mu_{\min}$. $\Pr(\mu(s, a) \leq b|\mu) = 0$. Thus
 266 we have:

267 **Corollary 3.** *With finite state action space and $b \leq \mu_{\min}$, under the assumptions as Theorem 1, $\hat{\pi}_t$*
 268 *from Algorithm 1 satisfies that with probability at least $1 - 3\delta$,*

$$v_M^{\mu} - v_M^{\hat{\pi}_t} \leq \frac{52V_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|} (\sqrt{\ln(2|\mathcal{S}||\mathcal{A}|/\delta)} + \sqrt{\ln(1+nV_{\max})}) + 8}{\sqrt{nb}(1-\gamma)^3} \\ + \frac{12V_{\max}|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{nb(1-\gamma)^3} + \frac{3\gamma^{K-1}V_{\max}}{(1-\gamma)^2}$$

269 *Proof.* In finite state action space, the number of all deterministic policies is less than $|\mathcal{A}|^{|\mathcal{S}|}$. Thus
 270 we have a policy class with $\epsilon_{\Pi} = 0$ and $|\Pi| \leq |\mathcal{A}|^{|\mathcal{S}|}$. Since the Q value is bounded in $[0, V_{\max}]$,
 271 we can construct a ϵ covering set \mathcal{F} of all value functions in $[0, V_{\max}]^{|\mathcal{S}||\mathcal{A}|}$ with $(\frac{V_{\max}}{\epsilon} + 1)^{|\mathcal{S}||\mathcal{A}|}$
 272 functions. Then $\epsilon_{\mathcal{F}} \leq \max_g \min_{f \in \mathcal{F}} \|f - g\|_{\mu, 2} \leq \max_g \min_{f \in \mathcal{F}} \|f - g\|_{\infty} \leq \epsilon$.

273 We can also bound ϵ_{μ} in finite state action space. For any fixed s, a , by Bernstein's inequality we
 274 have that with probability of $1 - \frac{\delta}{|\mathcal{S}||\mathcal{A}|}$:

$$|\hat{\mu}(s, a) - \mu(s, a)| = \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}(s^{(i)} = s, a^{(i)} = a) - \mathbb{E}[\mathbb{1}(s^{(i)} = s, a^{(i)} = a)] \right| \quad (98)$$

$$\leq \sqrt{\frac{2\mathbb{V}[\mathbb{1}(s^{(i)} = s, a^{(i)} = a)] \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n}} + \frac{4 \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} \quad (99)$$

$$= \sqrt{\frac{2\mu(s, a)(1 - \mu(s, a)) \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n}} + \frac{4 \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} \quad (100)$$

275 By taking summation of $|\hat{\mu}(s, a) - \mu(s, a)|$ and union bound over all (s, a) , we can bound the total
 276 variation bounds between $\hat{\mu}$ and μ , with probability at least $1 - \delta$,

$$\|\hat{\mu} - \mu\|_{TV} = \frac{1}{2} \sum_{s,a} |\hat{\mu}(s, a) - \mu(s, a)| \quad (101)$$

$$\leq \frac{1}{2} \sum_{s,a} \left(\sqrt{\frac{2\mu(s, a)(1 - \mu(s, a)) \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n}} + \frac{4 \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} \right) \quad (102)$$

$$= \frac{2|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} + \frac{1}{2} \sum_{s,a} \sqrt{\frac{2\mu(s, a)(1 - \mu(s, a)) \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n}} \quad (103)$$

$$\leq \frac{2|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} + \frac{1}{2} \sqrt{\sum_{s,a} \frac{2\mu(s, a) \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} \sum_{s,a} (1 - \mu(s, a))} \quad (104)$$

(Cauchy-Schwartz's inequality)

$$= \frac{2|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} + \frac{1}{2} \sqrt{\frac{2 \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} (|\mathcal{S}||\mathcal{A}| - 1)} \quad (104)$$

$$\leq \frac{2|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} + \sqrt{\frac{|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{2n}} \quad (105)$$

277 Now in a finite state action space we can construct the policy and Q function sets with $|\mathcal{F}| \leq$
 278 $(\frac{V_{\max}}{\epsilon} + 1)^{|\mathcal{S}||\mathcal{A}|}$, $|\Pi| \leq |\mathcal{A}|^{|\mathcal{S}|}$, $\epsilon_{\Pi} = 0$, $\epsilon_{\mathcal{F}} \leq \epsilon$, and bounded ϵ_{μ} . By plugging these terms into the
 279 result of Theorem 1, we have the following bound:

$$v_M^{\mu} - v_M^{\hat{\pi}_t} \leq \frac{4C}{(1 - \gamma)^3} \left(\sqrt{\frac{419V_{\max}^2 (|\mathcal{S}| \ln |\mathcal{A}| + |\mathcal{S}||\mathcal{A}| \ln(1 + V_{\max}/\epsilon) + \ln(1/\delta))}{3n}} + 2\sqrt{\epsilon} \right) \\ + \frac{6CV_{\max}}{(1 - \gamma)^3} \left(\sqrt{\frac{|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{2n}} + \frac{2|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} \right) + \frac{3\gamma^{K-1}V_{\max}}{(1 - \gamma)^2}, \quad (106)$$

280 for any chosen $\epsilon > 0$. So we can set that $\epsilon = 1/n$ to upper bound the the infimum of this upper
 281 bound.

$$v_M^{\mu} - v_M^{\hat{\pi}_t} \leq \frac{4C}{(1 - \gamma)^3} \left(\sqrt{\frac{419V_{\max}^2 (|\mathcal{S}| \ln |\mathcal{A}| + |\mathcal{S}||\mathcal{A}| \ln(1 + nV_{\max}) + \ln(1/\delta))}{3n}} + 2\sqrt{\frac{1}{n}} \right) \\ + \frac{6CV_{\max}}{(1 - \gamma)^3} \left(\sqrt{\frac{|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{2n}} + \frac{2|\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{n} \right) + \frac{3\gamma^{K-1}V_{\max}}{(1 - \gamma)^2} \quad (107)$$

282 Notice that in discrete space we have that $U \leq 1$. By replacing C with $1/b$ and simplify some terms,
 283 we have that:

$$v_M^{\mu} - v_M^{\hat{\pi}_t} \leq \sqrt{\frac{6704V_{\max}^2 |\mathcal{S}| (\ln(|\mathcal{A}|/\delta) + |\mathcal{A}| \ln(1 + nV_{\max}))}{3nb^2(1 - \gamma)^6}} + \frac{8}{b\sqrt{n}(1 - \gamma)^3} \\ + \sqrt{\frac{18V_{\max}^2 |\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{nb^2(1 - \gamma)^6}} + \frac{12V_{\max} |\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{nb(1 - \gamma)^3} + \frac{3\gamma^{K-1}V_{\max}}{(1 - \gamma)^2} \\ \leq \frac{52V_{\max} \sqrt{|\mathcal{S}||\mathcal{A}| (\sqrt{\ln(2|\mathcal{S}||\mathcal{A}|/\delta)} + \sqrt{\ln(1 + nV_{\max}))}}}{\sqrt{nb}(1 - \gamma)^3} + 8 \\ + \frac{12V_{\max} |\mathcal{S}||\mathcal{A}| \ln(2|\mathcal{S}||\mathcal{A}|/\delta)}{nb(1 - \gamma)^3} + \frac{3\gamma^{K-1}V_{\max}}{(1 - \gamma)^2}$$

284

□

285 D Proofs for Q Iteration Guarantees

286 In this section, we are going to prove our main result for the Q iteration algorithm, Algorithm 2.
 287 First we introduce a similar completeness assumption about the Bellman optimality operator:

288 **Assumption 5** (Completeness under \mathcal{T}_ζ). $\max_{f \in \mathcal{F}} \min_{g \in \mathcal{F}} \|g - \mathcal{T}_\zeta f\|_{2,\mu}^2 \leq \epsilon_{\mathcal{F}}$

289 We will first state our main theorem here and then give a proof sketch before we start the proof
 290 formally.

291 **Theorem 4.** *Given a MDP $M = \langle \mathcal{S}, \mathcal{A}, R, P, \gamma, p \rangle$, a dataset $D = \{(s, a, r, s')\}$ with n samples
 292 that is drawn i.i.d. from $\mu \times R \times P$, and a finite Q -function classes \mathcal{F} satisfying Assumption 5, $\hat{\pi}_t$
 293 from Algorithm 2 satisfies that with probability at least $1 - \delta$, $v^{\tilde{\pi}} - v^{\hat{\pi}_t} \leq$*

$$\frac{2C}{(1-\gamma)^2} \left(\sqrt{\frac{208V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n}} + 2\sqrt{\epsilon_{\mathcal{F}}} + V_{\max}\epsilon_{\mu} + \|Q^{\tilde{\pi}} - \mathcal{T}_\zeta Q^{\tilde{\pi}}\|_{2,\mu} \right) + \frac{(2\gamma^t + \epsilon_\zeta)V_{\max}}{1-\gamma}$$

294 for any policy $\tilde{\pi} \in \Pi_C^{all}$.

295 We will first give a proof sketch before we start the proof formally. The proof follows a similar
 296 structural as the policy iteration case. To prove Theorem 4 we first prove a similar version of Theo-
 297 rem 4 but the comparator policies are in strong ζ -constrained policy set (formally stated as Theorem
 298 5 later). Then we show an upper bound of $v_{M'}^{\pi} - v_{M'}^{\pi_t}$, where $\pi \in \Pi_{SC}^{all}$ and π_t is the output of
 299 algorithm (Theorem 5, will be formally stated later). Then we are going to show that for any policy
 300 π in the ζ -constrained policy set, after a projection Ξ it is in the strong ζ -constrained policy set and
 301 $v_M^{\pi} \leq v_{M'}^{\Xi(\pi)} + V_{\max}\epsilon_\zeta/(1-\gamma)$. Then we can provide the upper bound for $v_M^{\pi} - v_M^{\pi_t}$ for any π in
 302 ζ -constrained policy set (Theorem 4).

The proof sketch of Theorem 5 goes as follow. One key step to prove this error bound is to convert
 the performance difference between any policy $\tilde{\pi} \in \Pi_{SC}^{all}$ and π_t to a value function gap that is
 filtered by ζ :

$$v^{\tilde{\pi}} - v^{\pi_t} \leq \|\zeta(Q^{\tilde{\pi}} - f_t)\|_{1,\nu_1}/(1-\gamma),$$

303 where ν_1 is some admissible distribution over $\mathcal{S} \times \mathcal{A}$. The filter ζ allows the change of measure from
 304 ν_1 to μ without constraining the density ratio between an arbitrary distribution ν and μ . Instead for
 305 any s, a where ζ is one, by definition μ is lower bounded and the density ratio is bounded by C
 306 (details in Lemma 13).

307 The rest of the proof has a similar structure with the standard FQI analysis. In Lemma 15, we bound
 308 the norm $\|\zeta(Q^{\tilde{\pi}} - f_t)\|_{2,\nu_1}$ by $C\|(f_t - \mathcal{T}_\zeta f_t)\|_{2,\mu}/(1-\gamma)$ and one additional sub-optimality error
 309 $\|Q^{\tilde{\pi}} - \mathcal{T}_\zeta Q^{\tilde{\pi}}\|_{2,\mu}$. The additional sub-optimality error term comes from the fact that $\tilde{\pi}$ may not be
 310 an optimal policy since the optimal policy may not be a ζ -constrained policy. The last step to finish
 311 the proof is to bound the expected Bellman residual by concentration inequality. Lemma 16 shows
 312 how to bound that following a similar approach as [1]. Then the main theorem is proved by combine
 313 all those steps. After that we prove when we can bound the value gap with respect to optimal value
 314 in Corollary 4.

315 Now we start the proof. We are going to condition on the high probability bounds in Assumption 2
 316 holds when we proof the lemmas.

Lemma 13. *For $\pi_t = \Xi(\hat{\pi}_t)$ in Algorithm 2, for any policy $\tilde{\pi} \in \Pi_{SC}^{all}$ we have*

$$v^{\tilde{\pi}} - v^{\pi_t} \leq \sum_{h=0}^{\infty} \gamma^h \left(\|\zeta(Q^{\tilde{\pi}} - f_t)\|_{1,\eta_h^{\pi_t} \times \tilde{\pi}} + \|\zeta(Q^{\tilde{\pi}} - f_t)\|_{1,\eta_h^{\tilde{\pi}} \times \pi_t} \right).$$

317 *Proof.* Given a deterministic greedy policy $\hat{\pi}_t$, $\pi_t = \Xi(\hat{\pi}_t)$ is also a deterministic policy and
 318 $\pi_t(s)$ equals $\hat{\pi}_t(s)$ unless $\zeta(s, \hat{\pi}_t(s)) = 0$, where $\pi_t(s) = a_{\text{abs}}$. Notice $\hat{\pi}_t(s)$ is the maximizer
 319 of $\zeta(s, \cdot)f_t(s, \cdot)$. If $\zeta(s, \hat{\pi}_t(s)) = 0$ then $\zeta(s, a)f_t(s, a) = 0$ for all a . We have that $\pi_t(s)$ is also the

320 maximizer of $\zeta(s, \cdot) f_t(s, \cdot)$.

$$v^{\tilde{\pi}} - v^{\pi_t} = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^{\pi_t}} [Q^{\tilde{\pi}}(s, \tilde{\pi}) - Q^{\tilde{\pi}}(s, \pi_t)] \quad ([3, \text{Lemma 6.1}])$$

$$\leq \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^{\pi_t}} \left[\zeta(s, \tilde{\pi}) Q^{\tilde{\pi}}(s, \tilde{\pi}) - \zeta(s, \pi_t) Q^{\tilde{\pi}}(s, \pi_t) \right] \quad (108)$$

$$\leq \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim \eta_h^{\pi_t}} \left[\zeta(s, \tilde{\pi}) Q^{\tilde{\pi}}(s, \tilde{\pi}) - \zeta(s, \tilde{\pi}) f_t(s, \tilde{\pi}) + \zeta(s, \pi_t) f_t(s, \pi_t) - \zeta(s, \pi_t) Q^{\tilde{\pi}}(s, \pi_t) \right] \quad (109)$$

$$\leq \sum_{h=0}^{\infty} \gamma^h \left(\left\| \zeta(Q^{\tilde{\pi}} - f_t) \right\|_{1, \eta_h^{\pi_t} \times \tilde{\pi}} + \left\| \zeta(Q^{\tilde{\pi}} - f_t) \right\|_{1, \eta_h^{\pi_t} \times \pi_t} \right) \quad (110)$$

321 Equation (108) follows from the fact that for all s, a such that $\tilde{\pi}(a|s) > 0$, either $\zeta(s, a) = 1$, or
 322 $a = a_{\text{abs}}$. $a = a_{\text{abs}}$ indicates that $Q^{\tilde{\pi}}(s, a) = 0$. So for all s, a such that $\tilde{\pi}(a|s) > 0$, $Q^{\tilde{\pi}}(s, a) =$
 323 $\zeta(s, a) Q^{\tilde{\pi}}(s, a)$. The second part follows from that for any s, a , $Q^{\tilde{\pi}}(s, a) \geq \zeta(s, a) Q^{\tilde{\pi}}(s, a)$. Equa-
 324 tion (109) follows from the fact that $\pi_t(s)$ is the maximizer of $\zeta(s, \cdot) f_t(s, \cdot)$. \square

Lemma 14. For any two function $f_1, f_2 : \mathcal{S}' \times \mathcal{A}' \rightarrow \mathbb{R}^+$, define $\pi_{f_1, f_2}(s) = \arg \max_{a \in \mathcal{A}} |f_1(s, a) - f_2(s, a)|$. Then we have $\forall \nu : \mathcal{S}' \rightarrow \Delta(\mathcal{A}')$,

$$\left\| \max_{a \in \mathcal{A}} f_1 - \max_{a \in \mathcal{A}} f_2 \right\|_{1, P(\nu)} \leq \|f_1 - f_2\|_{1, P(\nu) \times \pi_{f_1, f_2}}.$$

Proof.

$$\begin{aligned} \left\| \max_{a \in \mathcal{A}} f_1 - \max_{a \in \mathcal{A}} f_2 \right\|_{1, P(\nu)} &= \mathbb{E}_{s \sim P(\nu)} \left| \max_{a \in \mathcal{A}} f_1(s, a) - \max_{a \in \mathcal{A}} f_2(s, a) \right| \\ &\leq \mathbb{E}_{s \sim P(\nu)} \max_{a \in \mathcal{A}} |f_1(s, a) - f_2(s, a)| \\ &= \mathbb{E}_{s \sim P(\nu), a \sim \pi_{f_1, f_2}} |f_1(s, a) - f_2(s, a)| \\ &= \|f_1 - f_2\|_{1, P(\nu) \times \pi_{f_1, f_2}}^2. \end{aligned}$$

325 \square

326 **Lemma 15.** For the data distribution μ and any admissible distribution ν over $\mathcal{S}' \times \mathcal{A}'$, $f, f' : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^+$ and any $\tilde{\pi} \in \Pi_{\mathcal{S}\mathcal{C}}^{\text{all}}$, we have

$$\begin{aligned} \left\| \zeta(f - Q^{\tilde{\pi}}) \right\|_{1, \nu} &\leq C \left(\|f - \mathcal{T}_{\zeta} f'\|_{2, \mu} + \left\| \mathcal{T}_{\zeta} Q^{\tilde{\pi}} - Q^{\tilde{\pi}} \right\|_{2, \mu} + V_{\max} \epsilon_{\mu} \right) \\ &\quad + \gamma \left\| \zeta(f' - Q^{\tilde{\pi}}) \right\|_{2, P(\nu) \times \pi_{\zeta f', \zeta Q^{\tilde{\pi}}}}. \end{aligned}$$

Proof.

$$\left\| \zeta \left(f - Q^{\tilde{\pi}} \right) \right\|_{1,\nu} \quad (111)$$

$$= \left\| \zeta \left(f - \mathcal{T}_\zeta f' + \mathcal{T}_\zeta f' - \mathcal{T}_\zeta Q^{\tilde{\pi}} + \mathcal{T}_\zeta Q^{\tilde{\pi}} - Q^{\tilde{\pi}} \right) \right\|_{1,\nu} \quad (112)$$

$$\leq \left\| \zeta \left(f - \mathcal{T}_\zeta f' \right) \right\|_{1,\nu} + \left\| \zeta \left(\mathcal{T}_\zeta f' - \mathcal{T}_\zeta Q^{\tilde{\pi}} \right) \right\|_{1,\nu} + \left\| \zeta \left(\mathcal{T}_\zeta Q^{\tilde{\pi}} - Q^{\tilde{\pi}} \right) \right\|_{1,\nu} \quad (113)$$

$$\leq C \left\| f - \mathcal{T}_\zeta f' \right\|_{1,\hat{\mu}} + \gamma \left\| \max_{a \in \mathcal{A}} \zeta f' - \max_{a \in \mathcal{A}} \zeta Q^{\tilde{\pi}} \right\|_{1,P(\nu)} + C \left\| \mathcal{T}_\zeta Q^{\tilde{\pi}} - Q^{\tilde{\pi}} \right\|_{1,\hat{\mu}} \quad (114)$$

$$\leq 2CV_{\max} \epsilon_\mu + C \left\| f - \mathcal{T}_\zeta f' \right\|_{1,\mu} + \gamma \left\| \max_{a \in \mathcal{A}} \zeta f' - \max_{a \in \mathcal{A}} \zeta Q^{\tilde{\pi}} \right\|_{1,P(\nu)} + C \left\| \mathcal{T}_\zeta Q^{\tilde{\pi}} - Q^{\tilde{\pi}} \right\|_{1,\mu} \quad (115)$$

$$\leq C \left(\left\| f - \mathcal{T}_\zeta f' \right\|_{2,\mu} + \left\| \mathcal{T}_\zeta Q^{\tilde{\pi}} - Q^{\tilde{\pi}} \right\|_{1,\mu} + 2V_{\max} \epsilon_\mu \right) + \gamma \left\| \zeta \left(f' - Q^{\tilde{\pi}} \right) \right\|_{1,P(\nu) \times \pi_{\zeta f', \zeta Q^{\tilde{\pi}}}} \quad (116)$$

328 The change of norms from $\| \cdot \|_\nu$ to $\| \cdot \|_\mu$ follows from that $\zeta(s, a) \neq 0$ iff $\hat{\mu}(s, a) \geq b$ and thus
 329 $\nu(s, a) \leq \hat{\mu}(s, a)U/b = C\hat{\mu}(s, a)$. The last step follows from Lemma 14. $\left\| \zeta \left(\mathcal{T}_\zeta f' - \mathcal{T}_\zeta Q^{\tilde{\pi}} \right) \right\|_{1,\nu} \leq$
 330 $\gamma \left\| \max_{a \in \mathcal{A}} \zeta f' - \max_{a \in \mathcal{A}} \zeta Q^{\tilde{\pi}} \right\|_{1,P(\nu)}$ follows from:

$$\left\| \zeta \left(\mathcal{T}_\zeta f' - \mathcal{T}_\zeta Q^{\tilde{\pi}} \right) \right\|_{1,\nu} = \mathbb{E}_{(s,a) \sim \nu} \left[\zeta(s, a) \left| \mathcal{T}_\zeta f'(s, a) - \mathcal{T}_\zeta Q^{\tilde{\pi}}(s, a) \right| \right] \quad (117)$$

$$\leq \mathbb{E}_{(s,a) \sim \nu} \left[\left| \mathcal{T}_\zeta f'(s, a) - \mathcal{T}_\zeta Q^{\tilde{\pi}}(s, a) \right| \right] \quad (118)$$

$$= \mathbb{E}_{(s,a) \sim \nu} \left[\left| \gamma \mathbb{E}_{s' \sim P(s,a)} \max_{a' \in \mathcal{A}} \zeta(s', a') f'(s', a') - \max_{a' \in \mathcal{A}} \zeta(s', a') Q^{\tilde{\pi}}(s', a') \right| \right] \quad (119)$$

$$\leq \gamma \mathbb{E}_{(s,a) \sim \nu, s' \sim P(s,a)} \left[\left| \max_{a' \in \mathcal{A}} \zeta(s', a') f'(s', a') - \max_{a' \in \mathcal{A}} \zeta(s', a') Q^{\tilde{\pi}}(s', a') \right| \right] \quad (Jensen)$$

$$= \gamma \mathbb{E}_{s' \sim P(\nu)} \left[\left| \max_{a' \in \mathcal{A}} \zeta(s', a') f'(s', a') - \max_{a' \in \mathcal{A}} \zeta(s', a') Q^{\tilde{\pi}}(s', a') \right| \right] \quad (120)$$

$$= \gamma \left\| \max_{a \in \mathcal{A}} \zeta f' - \max_{a \in \mathcal{A}} \zeta Q^{\tilde{\pi}} \right\|_{1,P(\nu)} \quad (121)$$

331

□

332 Now we are going to use Bernstein's inequality to bound $\|f_{t+1} - \mathcal{T}_\zeta f_t\|_{2,\mu}$, which mostly follows
 333 from [1]'s proof for the vanilla value iteration.

334 **Lemma 16.** *With Assumption 5 holds, let $g_f^* = \arg \min_{g \in \mathcal{F}} \|g - \mathcal{T}_\zeta f\|_{2,\mu}$, then $\|g_f^* - \mathcal{T}_\zeta f\|_{2,\mu}^2 \leq \epsilon_{\mathcal{F}}$.
 335 The dataset D is generated i.i.d. from M as follows: $(s, a) \sim \mu$, $r = R(s, a)$, $s' \sim P(s, a)$. Define
 336 $\mathcal{L}_\mu(f; f') = \mathbb{E}[\mathcal{L}_D(f; f')]$. We have that $\forall f \in \mathcal{F}$, with probability at least $1 - \delta$,*

$$\mathcal{L}_\mu(\mathcal{T}_{\zeta,D} f; f) - \mathcal{L}_\mu(g_f^*; f) \leq \frac{208V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n} + \epsilon_{\mathcal{F}}$$

337 where $\mathcal{T}_{\zeta,D} f = \arg \min_{g \in \mathcal{F}} \mathcal{L}_D(g, f)$.

338 *Proof.* This proof is similar with the proof of Lemma 7, and we adapt it to operator \mathcal{T}_ζ . The only
 339 change is the definition of $V_f(\cdot)$ and $X(\cdot, \cdot, \cdot)$. The definition of \mathcal{L}_D and \mathcal{L}_μ would not change
 340 between M and M' , and the right hand side is also the same constant for M and M' . So the result
 341 we prove here does not change from M to M' .

For the simplicity of notations, let $V_f(s) = \max_{a \in \mathcal{A}} \zeta(s, a) f(s, a)$. Fix $f, g \in \mathcal{F}$, and define

$$X(g, f, g_f^*) := (g(s, a) - r - \gamma V_f(s'))^2 - (g_f^*(s, a) - r - \gamma V_f(s'))^2.$$

Plugging each $(s, a, r, s') \in D$ into $X(g, f, g_f^*)$, we get i.i.d. variables $X_1(g, f, g_f^*), X_2(g, f, g_f^*), \dots, X_n(g, f, g_f^*)$. It is easy to see that

$$\frac{1}{n} \sum_{i=1}^n X_i(g, f, g_f^*) = \mathcal{L}_D(g; f) - \mathcal{L}_D(g_f^*; f).$$

By the definition of \mathcal{L}_μ , it is also easy to see that

$$\mathcal{L}_\mu(g; f) = \|g - \mathcal{T}_\zeta f\|_{2,\mu}^2 + \mathbb{E}_{s,a \sim \mu} \left[\mathbb{V}_{r,s'} \left(r + \gamma \max_{a' \in \mathcal{A}} \zeta(s', a') f(s', a') \right) \right]$$

Notice that the second part does not depend on g . Then

$$\mathcal{L}_\mu(g; f) - \mathcal{L}_\mu(\mathcal{T}_\zeta f; f) = \|g - \mathcal{T}_\zeta f\|_{2,\mu}^2$$

342 Then we bound the variance of X :

$$\begin{aligned} \mathbb{V}[X(g, f, g_f^*)] &\leq \mathbb{E}[X(g, f, g_f^*)^2] \\ &= \mathbb{E}_\mu \left[\left((g(s, a) - r - \gamma V_f(s'))^2 - (g_f^*(s, a) - r - \gamma V_f(s'))^2 \right)^2 \right] \\ &= \mathbb{E}_\mu \left[(g(s, a) - g_f^*(s, a))^2 (g(s, a) + g_f^*(s, a) - 2r - 2\gamma V_f(s'))^2 \right] \\ &\leq 4V_{\max}^2 \mathbb{E}_\mu \left[(g(s, a) - g_f^*(s, a))^2 \right] \\ &= 4V_{\max}^2 \|g - g_f^*\|_{2,\mu}^2 \tag{122} \\ &\leq 8V_{\max}^2 (\mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F}}). \tag{*} \end{aligned}$$

343 Step (*) holds because

$$\begin{aligned} &\|g - g_f^*\|_{2,\mu}^2 \\ &\leq 2 (\|g - \mathcal{T}_\zeta f\|_{2,\mu}^2 + \|\mathcal{T}_\zeta f - g_f^*\|_{2,\mu}^2) \tag{((a+b)^2 \leq 2a^2 + 2b^2)} \\ &\leq 2 (\|g - \mathcal{T}_\zeta f\|_{2,\mu}^2 - \|\mathcal{T}_\zeta f - g_f^*\|_{2,\mu}^2 + 2\|\mathcal{T}_\zeta f - g_f^*\|_{2,\mu}^2) \\ &= 2 [(\mathcal{L}_\mu(g; f) - \mathcal{L}_\mu(\mathcal{T}_\zeta f; f)) - (\mathcal{L}_\mu(g_f^*; f) - \mathcal{L}_\mu(\mathcal{T}_\zeta f; f)) + 2\|\mathcal{T}_\zeta f - g_f^*\|_{2,\mu}^2] \\ &= 2 [(\mathcal{L}_\mu(g; f) - \mathcal{L}_\mu(g_f^*; f)) + 2\|\mathcal{T}_\zeta f - g_f^*\|_{2,\mu}^2] \\ &= 2 (\mathbb{E}[X(g, f, g_f^*)] + 2\|\mathcal{T}_\zeta f - g_f^*\|_{2,\mu}^2) \\ &\leq 2(\mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F}}) \end{aligned}$$

344 Next, we apply (one-sided) Bernstein's inequality and union bound over all $f \in \mathcal{F}$ and $g \in \mathcal{F}$. With
345 probability at least $1 - \delta$, we have

$$\begin{aligned} \mathbb{E}[X(g, f, g_f^*)] - \frac{1}{n} \sum_{i=1}^n X_i(g, f, g_f^*) &\leq \sqrt{\frac{2\mathbb{V}[X(g, f, g_f^*)] \ln \frac{|\mathcal{F}|^2}{\delta}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}|^2}{\delta}}{3n} \\ &= \sqrt{\frac{32V_{\max}^2 (\mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F}}) \ln \frac{|\mathcal{F}|}{\delta}}{n}} + \frac{8V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n} \end{aligned}$$

Since $\mathcal{T}_{\zeta, D} f$ minimizes $\mathcal{L}_D(\cdot; f)$, it also minimizes $\frac{1}{n} \sum_{i=1}^n X_i(\cdot, f, g_f^*)$. This is because the two objectives only differ by a constant $\mathcal{L}_D(g_f^*; f)$. Hence,

$$\frac{1}{n} \sum_{i=1}^n X_i(\mathcal{T}_{\zeta, D} f, f, g_f^*) \leq \frac{1}{n} \sum_{i=1}^n X_i(g_f^*, f, g_f^*) = 0.$$

346 Then,

$$\mathbb{E}[X(\mathcal{T}_{\zeta, D} f, f, g_f^*)] \leq \sqrt{\frac{32V_{\max}^2 (\mathbb{E}[X(\mathcal{T}_{\zeta, D} f, f, g_f^*)] + 2\epsilon_{\mathcal{F}}) \ln \frac{|\mathcal{F}|}{\delta}}{n}} + \frac{8V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n}.$$

347 Solving for the quadratic formula,

$$\begin{aligned}
\mathbb{E}[X(\mathcal{T}_{\zeta, D} f, f, g_f^*)] &\leq \sqrt{48 \left(\frac{8V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n} \right)^2 + \frac{64V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{n} \epsilon_{\mathcal{F}} + \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n}} \\
&\leq \frac{(56 + 32\sqrt{3})V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n} + \sqrt{\frac{64V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{n} \epsilon_{\mathcal{F}}} \\
&\quad (\sqrt{a+b} \leq \sqrt{a} + \sqrt{b} \text{ and } \ln \frac{|\mathcal{F}|}{\delta} > 0) \\
&\leq \frac{112V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n} + \sqrt{\frac{64V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{n} \epsilon_{\mathcal{F}}} \\
&\leq \frac{112V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n} + \frac{32V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{n} + \epsilon_{\mathcal{F}} \\
&\leq \frac{208V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n} + \epsilon_{\mathcal{F}}
\end{aligned}$$

348 Noticing that $\mathbb{E}[X(\mathcal{T}_{\zeta, D} f; f, g_f^*)] = \mathcal{L}_{\mu}(\mathcal{T}_{\zeta, D} f; f) - \mathcal{L}_{\mu}(g_f^*; f)$, we complete the proof. \square

349 Now we could prove the main theorem about fitted Q iteration.

350 **Theorem 5.** Given a MDP $M = \langle \mathcal{S}, \mathcal{A}, R, P, \gamma, p \rangle$, a dataset $D = \{(s, a, r, s')\}$ with n samples
351 that is draw i.i.d. from $\mu \times R \times P$, and a finite Q-function classes \mathcal{F} satisfying Assumption 5,
352 $\pi_t = \Xi(\hat{\pi}_t)$ from Algorithm 2 satisfies that with probability at least $1 - 2\delta$, $v^{\tilde{\pi}} - v^{\pi_t} \leq$

$$\frac{2C}{(1-\gamma)^2} \left(\sqrt{\frac{208V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n}} + 2\sqrt{\epsilon_{\mathcal{F}}} + V_{\max}\epsilon_{\mu} + \|Q^{\tilde{\pi}} - \mathcal{T}_{\zeta}Q^{\tilde{\pi}}\|_{1, \mu} \right) + \frac{2\gamma^t V_{\max}}{1-\gamma}$$

353 for any policy $\tilde{\pi} \in \Pi_{SC}^{all}$.

Proof. Firstly, we can let $f = f_t$ and $f' = f_{t-1}$ in Lemma 15. This gives us that

$$\|f_t - Q^{\tilde{\pi}}\|_{1, \nu} \leq C \left(\|f_t - \mathcal{T}_{\zeta}f_{t-1}\|_{2, \mu} + \|Q^{\tilde{\pi}} - \mathcal{T}_{\zeta}Q^{\tilde{\pi}}\|_{1, \mu} + 2V_{\max}\epsilon_{\mu} \right) + \gamma \|f_{t-1} - Q^{\tilde{\pi}}\|_{1, P(\nu) \times \pi_{f_{t-1}, Q^{\tilde{\pi}}}}$$

354 Note that we can apply the same analysis on $P(\nu) \times \pi_{f_{t-1}, Q^{\tilde{\pi}}}$ and expand the inequality t times. It
355 then suffices to upper bound $\|f_t - \mathcal{T}_{\zeta}f_{t-1}\|_{2, \mu}$.

$$\begin{aligned}
&\|f_t - \mathcal{T}_{\zeta}f_{t-1}\|_{2, \mu}^2 \\
&= \mathcal{L}_{\mu}(f_t; f_{t-1}) - \mathcal{L}_{\mu}(\mathcal{T}_{\zeta}f_{t-1}; f_{t-1}) \quad (\text{Definition of } \mathcal{L}_{\mu}) \\
&= [\mathcal{L}_{\mu}(f_t; f_{t-1}) - \mathcal{L}_{\mu}(g_{f_{t-1}}^*; f_{t-1})] + [\mathcal{L}_{\mu}(g_{f_{t-1}}^*; f_{t-1}) - \mathcal{L}_{\mu}(\mathcal{T}_{\zeta}f_{t-1}; f_{t-1})] \\
&\leq \epsilon_4 + \|g_{f_{t-1}}^* - \mathcal{T}_{\zeta}f_{t-1}\|_{2, \mu}^2 \quad (\text{Lemma 16 and definition of } \mathcal{L}_{\mu}) \\
&\leq \epsilon_4 + \epsilon_{\mathcal{F}}. \quad (\text{Definition of } g_{Q_{k-1}}^* \text{ and Assumption 5})
\end{aligned}$$

The inequality holds with probability at least $1 - \delta$ and $\epsilon_4 = \frac{208V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n} + \epsilon_{\mathcal{F}}$. Noticing that ϵ_4 and $\epsilon_{\mathcal{F}}$ do not depend on t , and the inequality holds simultaneously for different t , we have that

$$\|f_t - Q^{\tilde{\pi}}\|_{1, \nu} \leq \frac{1-\gamma^t}{1-\gamma} C \left(\sqrt{(\epsilon_4 + \epsilon_{\mathcal{F}})} + V_{\max}\epsilon_{\mu} + \|Q^{\tilde{\pi}} - \mathcal{T}_{\zeta}Q^{\tilde{\pi}}\|_{1, \mu} \right) + \gamma^t V_{\max}.$$

356 Applying this to Lemma 13, we have that

$$\begin{aligned}
& v^{\tilde{\pi}} - v^{\pi_t} \\
& \leq \frac{2}{1-\gamma} \left(\frac{1-\gamma^t}{1-\gamma} C \left(\sqrt{(\epsilon_4 + \epsilon_{\mathcal{F}})} + V_{\max} \epsilon_{\mu} + \left\| Q^{\tilde{\pi}} - \mathcal{T}_{\zeta} Q^{\tilde{\pi}} \right\|_{1,\mu} \right) + \gamma^t V_{\max} \right) \\
& \leq \frac{2C}{(1-\gamma)^2} \left(\sqrt{\epsilon_4 + \epsilon_{\mathcal{F}}} + V_{\max} \epsilon_{\mu} + \left\| Q^{\tilde{\pi}} - \mathcal{T}_{\zeta} Q^{\tilde{\pi}} \right\|_{1,\mu} \right) + \frac{2\gamma^t V_{\max}}{1-\gamma} \\
& \leq \frac{2C}{(1-\gamma)^2} \left(\sqrt{\frac{208V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n}} + 2\sqrt{\epsilon_{\mathcal{F}}} + V_{\max} \epsilon_{\mu} + \left\| Q^{\tilde{\pi}} - \mathcal{T}_{\zeta} Q^{\tilde{\pi}} \right\|_{1,\mu} \right) + \frac{2\gamma^t V_{\max}}{1-\gamma}.
\end{aligned}$$

357

□

358 Now we are going to use the fact that there is an no-value-loss projection from the ζ -constrained
359 policy set to the strong ζ -constrained policy set to prove an error bound w.r.t any $\tilde{\pi} \in \Pi_C^{all}$.

360 **Theorem 2.** Given a MDP $M = \langle \mathcal{S}, \mathcal{A}, R, P, \gamma, p \rangle$, a dataset $D = \{(s, a, r, s')\}$ with n samples
361 that is draw i.i.d. from $\mu \times R \times P$, and a finite Q -function classes \mathcal{F} satisfying Assumption 5, $\hat{\pi}_t$
362 from Algorithm 2 satisfies that with probability at least $1 - 2\delta$, $v^{\tilde{\pi}} - v^{\hat{\pi}_t} \leq$

$$\frac{2C}{(1-\gamma)^2} \left(\sqrt{\frac{208V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n}} + 2\sqrt{\epsilon_{\mathcal{F}}} + V_{\max} \epsilon_{\mu} + \left\| Q^{\tilde{\pi}} - \mathcal{T}_{\zeta} Q^{\tilde{\pi}} \right\|_{2,\mu} \right) + \frac{(2\gamma^t + \epsilon_{\zeta})V_{\max}}{1-\gamma}$$

363 for any policy $\tilde{\pi} \in \Pi_C^{all}$.

364 *Proof.* The difference between this theorem and Theorem 5 is that $\tilde{\pi}$ is in Π_C^{all} which is significantly
365 larger than Π_{SC}^{all} .

366 This prove mimics the proof of Theorem 1. For any policy $\tilde{\pi} \in \Pi_C^{all}$, Lemma 3 tells that $v_M^{\tilde{\pi}} \leq$
367 $v_{M'}^{\Xi(\tilde{\pi})} + \frac{V_{\max} \epsilon_{\zeta}}{1-\gamma}$. Since $\pi_t = \Xi(\hat{\pi}_t)$, $v_M^{\hat{\pi}_t} = v_{M'}^{\hat{\pi}_t} \geq v_M^{\pi_t}$. Then $v_M^{\tilde{\pi}} - v_M^{\hat{\pi}_t} \leq v_{M'}^{\Xi(\tilde{\pi})} - v_{M'}^{\hat{\pi}_t} + \frac{V_{\max} \epsilon_{\zeta}}{1-\gamma}$ and
368 Theorem 5 completes the proof. □

369 **Remark:** The first term in the theorem comes from that the best policy in the ζ -constrained policy
370 set is not optimal. Note that the ζ -constrained policy set does not requires any realizability to do
371 with our function approximation but merely about the density ratio of a policy. When there is an
372 optimal policy of M such in Π_C^{all} , we have the same type of bound as standard approximate value
373 iteration analysis.

374 **Corollary 4.** If there exists an π^* on M such that $\Pr(\mu(s, a) \leq 2b|\pi^*) \leq \epsilon$. then under the condition
375 as Theorem 4, $\hat{\pi}_t$ from Algorithm 2 satisfies that with probability at least $1 - 2\delta$, $v_M^{\pi^*} - v_M^{\hat{\pi}_t} \leq$

$$\frac{2C}{(1-\gamma)^2} \left(\sqrt{\frac{208V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta}}{3n}} + 2\sqrt{\epsilon_{\mathcal{F}}} + V_{\max} \epsilon_{\mu} + \left\| Q^{\pi^*} - \mathcal{T}_{\zeta} Q^{\pi^*} \right\|_{2,\mu} \right) + \frac{V_{\max}(2\gamma^t + \epsilon + CU\epsilon_{\mu})}{1-\gamma}$$

376 *Proof.* The proof of $\pi^* \in \Pi_C^{all}$ is same as the proof in Corollary 1. Then proof is finished by
377 applying Theorem 4. □

378 E Details of CartPole Experiment

379 E.1 Full results of Discretized CartPole-v0

380 In section 6.1, we compare AVI, BCQL[2], SPIBB[4], Behavior cloning and our algorithm PQI, in
381 CartPole-v0 with discretized state space. The data is generated by a ϵ -greedy policy (ϵ from 0.1 to
382 0.9) and we report the resulting policies from different algorithm with the best hyper-parameter in
383 each ϵ . In this section we show the learning curve for each ϵ and each hyper-parameter value. We
384 run the BCQ algorithm with the threshold of $\hat{\mu}(a|s)$ in $\{0, 0.05, 0.1, 0.2\}$, and we run the SPIBB

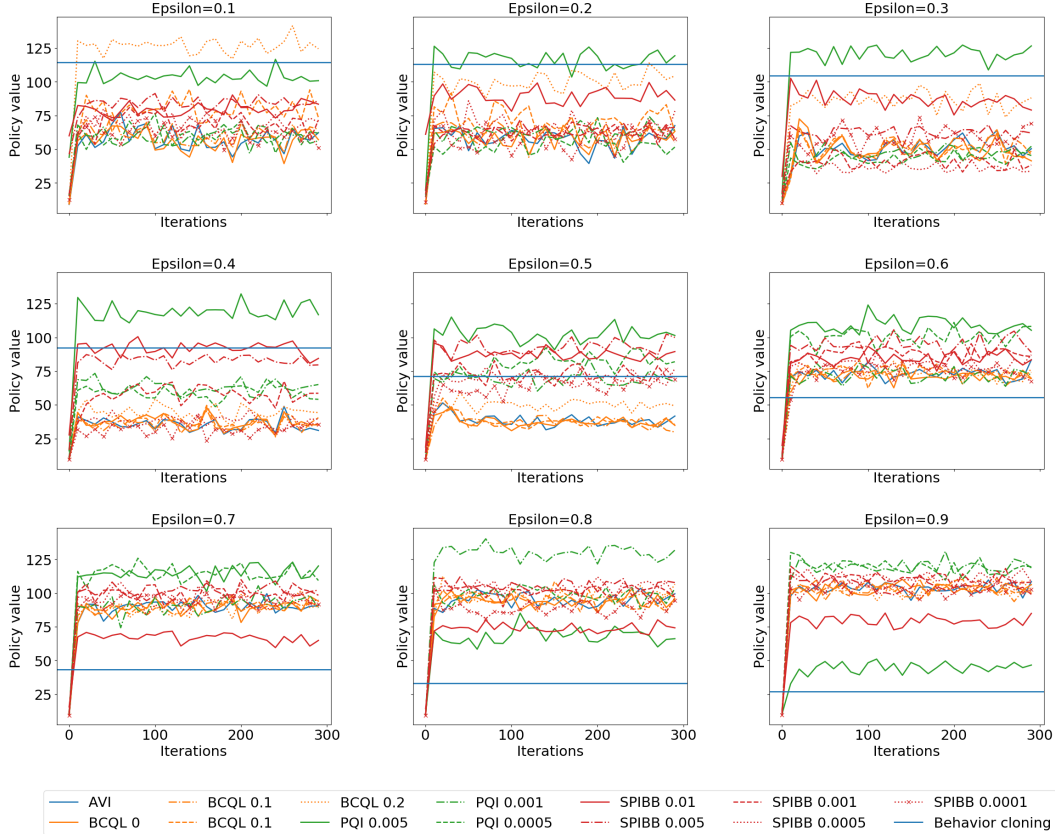


Figure 1: CartPole-v0 with discretized state space. The learning curve of all algorithms with different hyper-parameters, data generated with different ϵ -greedy behavior policy. The hyper-parameter of SPIBB [4] and PQI is the threshold of $\hat{\mu}(s, a)$ and the hyper-parameter of BCQL [2] is the threshold of $\hat{\mu}(a|s)$.

385 algorithm with the threshold of $\hat{\mu}(s, a)$ in $\{0.01, 0.005, 0.001, 0.0005, 0.0001\}$ and PQI with the
 386 threshold of $\hat{\mu}(s, a)$ in a smaller set $\{0.005, 0.001, 0.0005\}$. Figure 1 shows for most of the ϵ and
 387 threshold our algorithm tie with the best baseline (SPIBB), and the best threshold of our algorithm
 388 outperform all baseline algorithms in 8 out of 9 cases.

389 In Figure 1, we observe the trend that smaller ϵ will prefer a smaller b . This is verified by more
 390 results in the next section, and we discuss the reasons for this phenomenon there.

391 E.2 Ablation study of threshold b

392 A key aspect of our algorithm is to filter the state space by a threshold on the estimated probability
 393 $\hat{\mu}(s, a)$. This prevents the algorithm from updating using low-confidence state, action pairs when
 394 bootstrapping values. Then the choice of threshold b is a key trade-off in our algorithm: if b is too
 395 small it can not remove the low-confident state, action pairs effectively; if b is too large it might
 396 remove too many state, action pairs and prevent learning from more data. In order to demonstrate
 397 the effect of b and how should we choose b in different settings, we show the performance of PQI in
 398 a large range of b and several ϵ values.

399 In figure 2 we show the trend that smaller b works better for larger ϵ and larger b works better for
 400 smaller ϵ in general. This can be explained in the following way: with a larger ϵ the data distribution
 401 is more exploratory and hence the probabilities on individual state, action pairs are smaller. So a the
 402 same threshold that performs well with low exploration now censors a much larger part of the state,
 403 action space, necessitating a smaller threshold as ϵ is increased. In general, we find that having the
 404 largest threshold which still retains a significant fraction of the state, action space is a good heuristic
 405 for setting the b parameter.

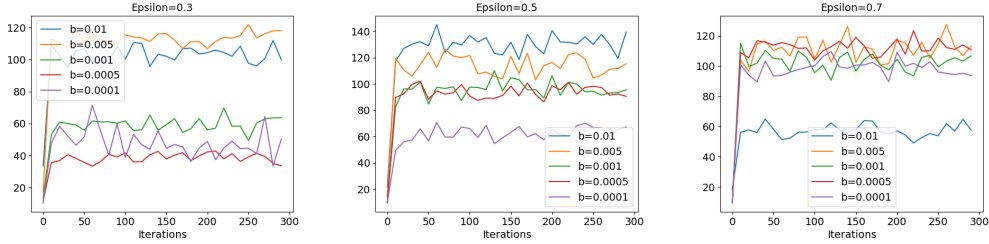


Figure 2: Performance of PQI with different values of threshold b

406 F Details of D4RL Experiment

407 In this section we introduce some missing details about the PQL algorithm and the experimental
 408 details in D4RL tasks. Our code is available at <https://github.com/yaoliucs/PQL>.

409 PQL algorithm is implemented based on the architecture of Batch-Constrained deep Q -learning
 410 (BCQ) [2] algorithm. More specifically, we use the similar Clipped Double Q-Learning (CDQ) up-
 411 date rule for the Q learning part, and employ a similar variational auto-encoder to fit the conditional
 412 action distribution in the batch. We use an additional variational auto-encoder to fit the marginalized
 413 state distribution of the batch. To implement an actual Q learning algorithm instead of an actor-critic
 414 algorithm, we did not sample from the actor in the Bellman backup but sample a larger batch from
 415 the fitted conditional action distribution. Algorithm 4 shows the pseudo-code of PQL to provide
 416 more details. We highlight the difference with the BCQ algorithm in red.

Algorithm 4 Pessimistic Q -learning (PQL)

Input: Batch D , ELBO threshold b , maximum perturbation Φ , target update rate τ , mini-batch size N , max number of iteration T . Number of actions k .

Initialize two Q network Q_{θ_1} and Q_{θ_2} , policy (perturbation) model: ξ_ϕ . ($\xi_\phi \in [-\Phi, \Phi]$), action VAE $G_{\omega_1}^a$ and state VAE $G_{\omega_2}^s$.

Pretrain $G_{\omega_2}^s$: $\omega_2 \leftarrow \arg \min_{\omega_2} ELBO(B; G_{\omega_2}^s)$.

for $t = 1$ **to** T **do**

 Sample a minibatch B with N samples from D .

$\omega_1 \leftarrow \arg \min_{\omega_1} ELBO(B; G_{\omega_1}^a)$.

Sample k actions a'_i from $G_{\omega_1}^a(s')$ for each s' .

 Compute the target y for each (s, a, r, s') pair:

$$y = r + \gamma \mathbf{1}(ELBO(s'; G_{\omega_2}^s) \geq b) \left[\max_{a'_i} \left(0.75 * \min_{j=1,2} Q_{\theta_j} + 0.25 * \max_{j=1,2} Q_{\theta_j} \right) \right]$$

$\theta \leftarrow \arg \min_{\theta} \sum (y - Q_{\theta}(s, a))^2$

 Sample k actions a_i from $G_{\omega_1}^a(s)$ for each s .

$\phi \leftarrow \arg \max_{\phi} \sum \max_{a_i} Q_{\theta_1}(s, a_i + \xi_\phi(s, a_i))$

 Update target network: $\theta' = (1 - \tau)\theta' + \tau\theta$, $\phi' = (1 - \tau)\phi' + \tau\phi$

end for

When evaluate the resulting policy: select action $a = \arg \max_{a_i} Q_{\theta_1}(s, a_i + \xi_\phi(s, a_i))$ where a_i are k actions sampled from $G_{\omega_1}^a(s)$ given s .

417 In practice, the indicator function $\mathbf{1}(ELBO(s'; G_{\omega_2}^s) \geq b)$ is implemented by
 418 $\text{sigmoid}(100(ELBO(s'; G_{\omega_2}^s) - b))$ to provide a slightly more smooth target. The evidence
 419 lower bound (ELBO) in VAE is:

$$ELBO(s; G_{\omega_2}^s) = \sum (s - \tilde{s})^2 + D_{\text{KL}}(N(\mu, \sigma) || N(0, 1)) \quad (123)$$

420 where μ and σ is sampled from the encoder of VAE with input s and \tilde{s} is sampled from the de-
 421 coder with the hidden state generated from $N(\mu, \sigma)$. $ELBO(B; G_{\omega_2}^s)$ is the averaged ELBO on the
 422 minibatch B . So does $G_{\omega_1}^a$. Note that this ELBO objective make the implicit assumption that the
 423 decoder's distribution is a Gaussian distribution with mean equals to the output of decoder network.

424 So when we generate the sample a' for computing y , we add a Gaussian noise to recover a sample
425 from the full posterior distribution.

426 For most of the hyper-parameters in Algorithm 4, we use the same value with the BCQ algorithm.
427 We run all algorithms with $T = 5 \times 10^5$ gradient steps as other reported results in D4RL tasks, and
428 the minibatch size $N = 100$ at each step. The number of sampled action when running the policy is
429 $k = 100$. Target network update rate is 0.005. The threshold b of ELBO is selected as 2-percentile
430 of the $ELBO(s)$ in the whole dataset after pretrain the VAE.

431 **References**

- 432 [1] Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learn-
433 ing. In *ICML*, pages 1042–1051, 2019.
- 434 [2] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning with-
435 out exploration. In *ICML*, pages 2052–2062, 2019.
- 436 [3] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning.
437 In *ICML*, pages 267–274, 2002.
- 438 [4] Romain Laroche, Paul Trichelair, and Remi Tachet Des Combes. Safe policy improvement with
439 baseline bootstrapping. In *ICML*, pages 3652–3661, 2019.