
Bellman Residual Orthogonalization for Offline Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose and analyze a reinforcement learning principle that approximates the
2 Bellman equations by enforcing their validity only along an user-defined space of
3 test functions. Focusing on applications to model-free offline RL with function
4 approximation, we exploit this principle to derive confidence intervals for off-policy
5 evaluation, as well as to optimize over policies within a prescribed policy class.
6 We prove an oracle inequality on our policy optimization procedure in terms of
7 a trade-off between the value and uncertainty of an arbitrary comparator policy.
8 Different choices of test function spaces allow us to tackle different problems
9 within a common framework. We characterize the loss of efficiency in moving
10 from on-policy to off-policy data using our procedures, and establish connections
11 to concentrability coefficients studied in past work. We examine in depth the
12 implementation of our methods with linear function approximation, and provide
13 theoretical guarantees with polynomial-time implementations even when Bellman
14 closure does not hold.

15 1 Introduction

16 Markov decision processes (MDP) provide a general framework for optimal decision-making in
17 sequential settings (e.g., [Put94, Ber95a, Ber95b]). Reinforcement learning refers to a general
18 class of procedures for estimating near-optimal policies based on data from an unknown MDP
19 (e.g., [BT96, SB18]). Different classes of problems can be distinguished depending on our access
20 to the data-generating mechanism. Many modern applications of RL involve learning based on a
21 pre-collected or offline dataset. Moreover, the state-action spaces are often sufficiently complex that
22 it becomes necessary to implement function approximation. In this paper, we focus on model-free
23 offline reinforcement learning (RL) with function approximation, where prior knowledge about the
24 MDP is encoded via the value function. In this setting, we focus on two fundamental problems: (1)
25 offline policy evaluation—namely, the task of accurately predicting the value of a target policy; and
26 (2) offline policy optimization, which is the task of finding a high-performance policy.

27 There are various broad classes of approaches to off-policy evaluation, including importance sam-
28 pling [Pre00, TB16, JL16, LLTZ18], as well as regression-based methods [LP03, MS08, CJ19].
29 Many methods for offline policy optimization build on these techniques, with a line of recent pa-
30 pers including the addition of pessimism [JYW21, XCJ⁺21, ZWB21]. We provide a more detailed
31 summary of the literature in Appendix A.3.

32 In contrast, this work investigates a different model-free principle—different from importance
33 sampling or regression-based methods—to learn from an offline dataset. It belongs to the class
34 of weight learning algorithms, which leverage an auxiliary function class to either encode the
35 marginalized importance weights of the target policy [LLTZ18, XJ20b], or estimates of the Bellman
36 errors [ASM08, CJ19, XJ20b]. Some work has considered kernel classes [FRTL20] or other weight

37 classes to construct off-policy estimators [UHJ20] as well as confidence intervals at the population
38 level [JH20]. However, these works do not examine in depth the statistical aspects of the problem, nor
39 elaborate upon the design of the weight function classes.¹ The last two considerations are essential to
40 obtaining data-dependent procedures accompanied by rigorous guarantees, and to provide guidance
41 on the choice of weight class, which are key contributions of this paper.

42 For space reasons, we motivate our approach in the idealized case where the Bellman operator is
43 known in Appendix A.1, and compare with the weight learning literature at the population level
44 in Appendix A.2. Let us summarize our main contributions in the following three paragraphs.

45 **Conceptual contributions:** Our paper makes two novel contributions of conceptual nature:

- 46 1. We propose a method, based on *approximate empirical orthogonalization* of the Bellman residual
47 along test functions, to construct confidence intervals and to perform policy optimization.
- 48 2. We propose a sample-based approximation of such principle, based on *self-normalization* and
49 *regularization*, and obtain general guarantees for parametric as well as non-parametric problems.

50 The construction of the estimator, its statistical analysis, and the concrete consequences (described
51 in the next paragraph) are the major distinctions with respect to past work on weight learning
52 methods [UHJ20, JH20]. Our analysis highlights the statistical trade-offs in the choice of the test
53 functions. (See Appendix A.2 for comparison with past work at the population level.)

54 **Domain-specific results:** In order to illustrate the broad effectiveness and applicability of our general
55 method and analysis, we consider several domains of interest. We show how to recover various results
56 from past work—and to obtain novel ones—by making appropriate choices of the test functions and
57 invoking our main result. Among these consequences, we discuss the following:

- 58 1. When marginalized importance weights are available, they can be used as test class. In this case
59 we recover a similar results as the paper [XJ20b]; however, here we only require concentrability
60 with respect to a comparator policy instead of over all policies in the class.
- 61 2. When some knowledge of the Bellman error class is available, it can be used as test class. Similar
62 results have appeared previously either with stronger concentrability [CJ19] or in the special case
63 of Bellman closure [XCJ⁺21].
- 64 3. We provide a test class that projects the Bellman residual along the error space of the Q class. The
65 resulting procedure is as an extension of the LSTD algorithm [BB96] to non-linear spaces, which
66 makes it a natural approach if no domain-specific knowledge is available. A related result is the
67 lower bound by [FKSLX21], which proves that without Bellman closure learning is hard even
68 with small density ratios. In contrast, our work shows that learning is still possible even with large
69 density ratios.
- 70 4. Finally, our procedure inherits some form of “multiple robustness”. For example, the two test
71 classes corresponding to Bellman completeness and marginalized importance weights can be used
72 together, and guarantees will be obtained if *either* Bellman completeness holds or the importance
73 weights are correct. We examine this issue in Section 4.4.

74 **Linear setting:** We examine in depth an application to the linear setting, where we propose the first
75 *computationally tractable* policy optimization procedure *without assuming Bellman completeness*.
76 The closest result here is given in the paper [ZWB21], which holds under Bellman closure. Our
77 procedure can be thought of making use of LSTD-type estimates so as to establish confidence intervals
78 for the projected Bellman equations, and then using an iterative scheme for policy improvement.

79 2 Background and set-up

80 We begin with some notation used throughout the paper. For a given probability distribution ρ
81 over a space \mathcal{X} , we define the $L^2(\rho)$ -inner product and semi-norm as $\langle f_1, f_2 \rangle_\rho = \mathbb{E}_\rho[f_1 f_2]$, and
82 $\|f_1\|_\rho = \sqrt{\langle f_1, f_1 \rangle_\rho}$. The identity function that returns one for every input is denoted by $\mathbb{1}$. We
83 frequently use notation such as $c, c', \tilde{c}, c_1, c_2$ etc. to denote constants that can take on different values
84 in different sections of the paper.

¹For instance, the paper [FRTL20] only shows validity of their intervals, not a performance bound; on the other hand, the paper [JH20] gives analyses at the population level, and so does not address the alignment of weight functions with respect to the dataset in the construction of the empirical estimator, which we do via self-normalization and regularization. This precludes obtaining the same type of guarantees that we present here.

85 **2.1 Markov decision processes and Bellman errors**

86 We focus on infinite-horizon discounted Markov decision processes [Put94, BT96, SB18] with
 87 discount factor $\gamma \in [0, 1)$, state space \mathcal{S} , and an action set \mathcal{A} . For each state-action pair (s, a) , there
 88 is a reward distribution $R(s, a)$ supported in $[0, 1]$ with mean $r(s, a)$, and a transition $\mathbb{P}(\cdot \mid s, a)$.

89 A (stationary) stochastic policy π maps states to actions. For a given policy, its Q -function is the
 90 discounted sum of future rewards based on starting from the pair (s, a) , and then following the
 91 policy π in all future time steps $Q^\pi(s, a) = r(s, a) + \sum_{h=0}^{\infty} \gamma^h \mathbb{E}[r_h(S_h, A_h) \mid (S_0, A_0) = (s, a)]$,
 92 where the expectation is taken over trajectories with $A_h \sim \pi(\cdot \mid S_h)$, and $S_{h+1} \sim \mathbb{P}(\cdot \mid$
 93 $S_h, A_h)$ for $h = 1, 2, \dots$. We also use $Q^\pi(s, \pi) = \mathbb{E}_{A \sim \pi(\cdot \mid s)} Q^\pi(s, A)$ and define the *Bellman*
 94 *evaluation operator* as $(\mathcal{T}^\pi Q)(s, a) = r(s, a) + \mathbb{E}_{S^+ \sim \mathbb{P}(\cdot \mid s, a)} Q(S^+, \pi)$. The value function satisfies
 95 $V^\pi(s) = Q^\pi(s, \pi)$. In our analysis, we assume that policies have action-value functions that satisfy
 96 the uniform bound $\sup_{(s, a)} |Q^\pi(s, a)| \leq 1$. We are also interested in approximating optimal poli-
 97 cies, whose value and action-value functions are defined as $V^*(s) = V^{\pi^*}(s) = \sup_{\pi} V^\pi(s)$ and
 98 $Q^*(s, a) = Q^{\pi^*}(s, a) = \sup_{\pi} Q^\pi(s, a)$.

99 We assume that the starting state S_0 is drawn according to ν_{start} and study $V^\pi = \mathbb{E}_{S_0 \sim \nu_{\text{start}}}[V^\pi(S_0)]$.
 100 The *discounted occupancy measure* associated with a policy π is defined as $d_\pi(s, a) = (1 -$
 101 $\gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h[(S_h, A_h) = (s, a)]$. We adopt the shorthand notation \mathbb{E}_π for expectations over
 102 d_π . For any functions $f, g : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we make frequent use of the shorthands $\mathbb{E}_\pi[f] \stackrel{\text{def}}{=}$
 103 $\mathbb{E}_{(S, A) \sim d_\pi}[f(S, A)]$, and $\langle f, g \rangle_\pi \stackrel{\text{def}}{=} \mathbb{E}_{(S, A) \sim d_\pi}[f(S, A)g(S, A)]$. Note moreover that we have
 104 $\langle \mathbb{1}, f \rangle_\pi = \mathbb{E}_\pi[f]$ where $\mathbb{1}$ denotes the identity function.

105 For a given Q -function and policy π , let us define the *temporal difference error* (or TD error)
 106 associated with the sample $z = (s, a, r, s^+)$ and the *Bellman error* at (s, a)

$$(\delta^\pi Q)(z) \stackrel{\text{def}}{=} Q(s, a) - r - \gamma Q(s^+, \pi), \quad (\mathcal{B}^\pi Q)(s, a) \stackrel{\text{def}}{=} Q(s, a) - r(s, a) - \gamma \mathbb{E}_{S^+ \sim \mathbb{P}(s, a)} Q(s^+, \pi).$$

107 The TD error is a random variable function of z , while the Bellman error is its conditional expectation
 108 with respect to the immediate reward and successor state at (s, a) . Many of our bounds involve the
 109 quantity $\mathbb{E}_\pi \mathcal{B}^\pi Q = \mathbb{E}_{(S, A) \sim d_\pi} [\mathcal{B}^\pi Q(S, A)]$.

110 **2.2 Function Spaces and Weak Representation**

111 Our methods involve three different types of function spaces, corresponding to policies, action-
 112 value functions, and test functions. A test function f is a mapping $(s, a, o) \mapsto f(s, a, o)$ such
 113 that $\sup_{(s, a, o)} |f(s, a, o)| \leq 1$, where o is an optional identifier containing side information. Our
 114 methodology involves the following three function classes:

- 115 • a *policy class* Π that contains all policies π of interest (for evaluation or optimization);
- 116 • for each π , the *predictor class* \mathcal{Q}^π of action-value functions Q that we permit; and
- 117 • for each π , the *test function class* \mathcal{F}^π that we use to enforce the Bellman residual constraints.

118 We use the shorthands $\mathcal{Q} = \cup_{\pi \in \Pi} \mathcal{Q}^\pi$ and $\mathcal{F} = \cup_{\pi \in \Pi} \mathcal{F}^\pi$. We assume *weak realizability*:

119 **Assumption 1** (Weak Realizability). *For a given policy π , the predictor class \mathcal{Q}^π is weakly realizable*
 120 *with respect to the test space \mathcal{F}^π and the measure μ if there exists a predictor $Q_\star^\pi \in \mathcal{Q}^\pi$ such that*

$$\langle f, \mathcal{B}^\pi Q_\star^\pi \rangle_\mu = 0 \text{ for all } f \in \mathcal{F}^\pi \quad \text{and} \quad \langle \mathbb{1}, \mathcal{B}^\pi Q_\star^\pi \rangle_\pi = 0. \quad (1)$$

121 The first condition requires the predictor to satisfy the Bellman equations *on average*. The second
 122 condition amounts to requiring that the predictor returns the value of π at the start distribution: using
 123 Lemma 9 stated in the sequel, we have

$$\mathbb{E}_{S \sim \nu_{\text{start}}} Q_\star^\pi(S, \pi) - V^\pi = \mathbb{E}_{S \sim \nu_{\text{start}}} [Q_\star^\pi - Q^\pi](S, \pi) = \frac{1}{1 - \gamma} \mathbb{E}_\pi \mathcal{B}^\pi Q_\star^\pi = \frac{1}{1 - \gamma} \langle \mathbb{1}, \mathcal{B}^\pi Q_\star^\pi \rangle_\pi = 0.$$

124 This weak notion should be contrasted with *strong realizability*, which requires a function $Q^\pi \in \mathcal{Q}^\pi$
 125 that satisfies the Bellman equation in all state-action pairs.

126 A stronger assumption that we sometime use is Bellman closure, which requires that $\mathcal{T}^\pi(Q) \in$
 127 \mathcal{Q}^π for all $Q \in \mathcal{Q}^\pi$. The corresponding ‘weak’ version is given in Appendix A.4.

128 3 Policy Estimates via the Weak Bellman Equations

129 In this section, we introduce our high-level approach, first at the population level and then in terms of
130 regularized/normalized sample-based approximations.

131 3.1 Weak Bellman equations, empirical approximations and confidence intervals

132 We begin by noting that the predictor Q^π satisfies the Bellman equations everywhere in the state-
133 action space, i.e., $\mathcal{B}^\pi Q^\pi = 0$. However, if our dataset is “small” relative to the complexity of
134 (functions) on the state-action space, then it is unrealistic to enforce such a stringent condition.
135 Instead, the idea is to control the Bellman error in a weighted-average sense, where the weights are
136 given by a set of *test functions*. At the idealized population level (corresponding to an infinite sample
137 size), we consider predictors that satisfy the conditions

$$\langle f, \mathcal{B}^\pi Q \rangle_\mu = 0, \quad \text{for all } f \in \mathcal{F}^\pi. \quad (2)$$

138 where \mathcal{F}^π is a user-defined set of test functions. The two main challenges here are how to use data to
139 enforce an approximate version of such constraints (along with rigorous data-dependent guarantees),
140 and how to design the test function space. We begin with the former challenge.

141 **Construction of the empirical set:** Given a dataset $\mathcal{D} = \{(s_i, a_i, r_i, s_i^+, o_i)\}_{i=1}^n$, we can approxi-
142 mate the Bellman errors by a linear combination of the temporal difference errors:

$$\int \underbrace{f(s, a) [Q(s, a) - (\mathcal{T}^\pi Q)(s, a)]}_{=\mathcal{B}^\pi Q(s, a)} d\mu \approx \frac{1}{n} \sum_{i=1}^n \underbrace{f(s_i, a_i) [Q(s_i, a_i) - r_i - \gamma Q(s_i^+, \pi)]}_{=\delta^\pi Q(s_i, a_i, r_i, s_i^+, o_i)}. \quad (3)$$

143 Note that the approximation (3) corresponds to a weighted linear combination of temporal differences.
144 Written more compactly in inner product notation, equation (3) reads $\langle f, \mathcal{B}^\pi Q \rangle_\mu \approx \langle f, \delta^\pi Q \rangle_n$, where
145 $\langle f, g \rangle_n = \frac{1}{n} \sum_{(s, a, r, s^+, o) \in \mathcal{D}} (fg)(s, a, r, s^+, o)$.

146 In general, the action value function Q^π does not satisfy $\langle f, \delta^\pi Q^\pi \rangle_n = 0$ because the empirical
147 approximation (3) involves sampling error. For these reasons, in order to (approximately) identify
148 Q^π , we impose only inequalities. Given a class of test functions \mathcal{F}^π , a radius parameter $\rho \geq 0$ and
149 regularization parameter $\lambda \geq 0$, we define the set

$$\widehat{\mathcal{C}}_n^\pi(\rho, \lambda; \mathcal{F}^\pi) \stackrel{\text{def}}{=} \left\{ Q \in \mathcal{Q}^\pi \quad \text{such that} \quad \frac{|\langle f, \delta^\pi Q \rangle_n|}{\sqrt{\|f\|_n^2 + \lambda}} \leq \sqrt{\frac{\rho}{n}} \quad \text{for all } f \in \mathcal{F}^\pi \right\}. \quad (4)$$

150 When the choices of (ρ, λ) are clear from the context, we adopt the shorthand $\widehat{\mathcal{C}}_n^\pi(\mathcal{F}^\pi)$, or $\widehat{\mathcal{C}}_n^\pi$ when
151 the function class \mathcal{F}^π is also clear. If \mathcal{F}^π and \mathcal{Q}^π have finite cardinality, $\rho \approx \ln |\mathcal{F}^\pi| |\mathcal{Q}^\pi| + \ln 1/\delta$,
152 where δ is a prescribed failure probability.

153 Our definition of the empirical constraint set (4) has two key components: first, the division by
154 $\sqrt{\|f\|_n^2 + \lambda}$ corresponds to a form of *self-normalization*, whereas the addition of λ corresponds to a
155 form of *regularization*. Self-normalization is needed so that the constraints remain suitably scale-
156 invariant. More importantly—in conjunction with the regularization—it ensures that test functions
157 that have poor coverage under the dataset do not have major effects on the solution. In particular,
158 the empirical norm $\|f\|_n^2$ in the self-normalization measures how well the given test function is
159 covered by the dataset. Any test function with poor coverage (i.e., $\|f\|_n^2 \approx 0$) will not yield useful
160 information, and the regularization counteracts its influence. In our guarantees, the choices of λ and
161 ρ are critical; as shown in our theory, we typically have $\lambda = \rho/n$, where ρ scales with the metric
162 entropy of the predictor, test and policy spaces. Disregarding ρ , the right-hand side of the constraint
163 decays as $1/\sqrt{n}$, so that the constraints are enforced more tightly as the sample size increases.

164 **Confidence bounds and policy optimization:** First, for any fixed policy π , we can use the feasi-
165 bility set (4) to compute the lower and upper estimates

$$\widehat{V}_{\min}^\pi \stackrel{\text{def}}{=} \min_{Q \in \widehat{\mathcal{C}}_n^\pi(\rho, \lambda; \mathcal{F}^\pi)} \mathbb{E}_{S \sim \nu_{\text{start}}} [Q(S, \pi)], \quad \text{and} \quad \widehat{V}_{\max}^\pi \stackrel{\text{def}}{=} \max_{Q \in \widehat{\mathcal{C}}_n^\pi(\rho, \lambda; \mathcal{F}^\pi)} \mathbb{E}_{S \sim \nu_{\text{start}}} [Q(S, \pi)], \quad (5)$$

166 corresponding to estimates of the minimum and maximum value that the policy π can take at the
167 initial distribution. The family of lower estimates can be used to perform policy optimization over the
168 class Π , in particular by solving the *max-min* problem

$$\max_{\pi \in \Pi} \left[\min_{Q \in \widehat{\mathcal{C}}_n^\pi} \mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) \right], \quad \text{or equivalently} \quad \max_{\pi \in \Pi} \widehat{V}_{\min}^\pi. \quad (6)$$

169 **Form of guarantees** Let us now specify and discuss the types of guarantees that we establish for
 170 our estimators (5) and (6). All of our theoretical guarantees involve a μ -based counterpart \mathcal{C}_n^π of the
 171 data-dependent set $\widehat{\mathcal{C}}_n^\pi$. More precisely, we define the population set

$$\mathcal{C}_n^\pi(4\rho, \lambda; \mathcal{F}^\pi) \stackrel{def}{=} \left\{ Q \in \mathcal{Q}^\pi \text{ such that } \frac{|\langle f, \mathcal{B}^\pi Q \rangle_\mu|}{\sqrt{\|f\|_\mu^2 + \lambda}} \leq \sqrt{\frac{4\rho}{n}} \text{ for all } f \in \mathcal{F} \right\}, \quad (7)$$

172 where $\langle f, g \rangle_\mu \stackrel{def}{=} \int f(s, a)g(s, a)d\mu$ is the inner product induced by a distribution² μ over (s, a) .
 173 As before, we use the shorthand notation \mathcal{C}_n^π when the underlying arguments are clear from context.
 174 Moreover, in the sequel, we generally ignore the constant 4 in the definition (7) by assuming that ρ is
 175 rescaled appropriately—e.g., that we use a factor of $\frac{1}{4}$ in defining the empirical set.

176 It should be noted that in contrast to the set $\widehat{\mathcal{C}}_n^\pi$, the set \mathcal{C}_n^π is *non-random* and it is defined in terms of
 177 the distribution μ and the input space $(\Pi, \mathcal{F}, \mathcal{Q})$. It relaxes the orthogonality constraints in the weak
 178 Bellman formulation (2). Our guarantees for off-policy confidence intervals take the following form:

$$\text{Coverage guarantee: } \quad [\widehat{V}_{\min}^\pi, \widehat{V}_{\max}^\pi] \ni V^\pi. \quad (8a)$$

$$\text{Width bound: } \quad \max \left\{ |\widehat{V}_{\min}^\pi - V^\pi|, |\widehat{V}_{\max}^\pi - V^\pi| \right\} \leq \frac{1}{1-\gamma} \max_{Q \in \mathcal{C}_n^\pi(\mathcal{F}^\pi)} |\mathbb{E}_\pi \mathcal{B}^\pi Q|. \quad (8b)$$

179 Turning to policy optimization, let $\tilde{\pi}$ be a solution to the max-min criterion (6). Then we prove a
 180 result of the following type:

$$\text{Oracle inequality: } \quad V^{\tilde{\pi}} \geq \underbrace{\max_{\pi \in \Pi} \left\{ V^\pi \right\}}_{\text{Value}} - \underbrace{\frac{1}{1-\gamma} \max_{Q \in \mathcal{C}_n^\pi(\mathcal{F}^\pi)} |\mathbb{E}_\pi \mathcal{B}^\pi Q|}_{\text{Evaluation uncertainty}}. \quad (9)$$

181 Note that this result guarantees that the estimator competes against an oracle that can search over all
 182 policies, and select one based on the optimal trade-off between its value and evaluation uncertainty.

183 3.2 High-probability guarantees

184 In this section, we present some high-probability guarantees. So as to facilitate understanding under
 185 space constraints, we state here results under simplifying assumptions: (a) the dataset originates
 186 from a fixed distribution, and (b) the classes Π, \mathcal{F} and \mathcal{Q} have finite cardinality. We emphasize
 187 that Appendix B provides a far more general version of this result, with an extremely flexible
 188 sampling model, and involving metric entropies of parametric or non-parametric function classes.

189 **Assumption 2** (I.i.d. dataset). *An i.i.d. dataset is a collection $\mathcal{D} = \{(s_i, a_i, r_i, s_i^+, o_i)\}_{i=1}^n$ such that
 190 for each $i = 1, \dots, n$ we have $(s_i, a_i, o_i) \sim \mu$ and conditioned on (s_i, a_i, o_i) , we observe a noisy
 191 reward $r_i = r(s_i, a_i) + \eta_i$ with $\mathbb{E}[\eta_i | \mathcal{F}_i] = 0$, $|r_i| \leq 1$ and the next state $s_i^+ \sim \mathbb{P}(s_i, a_i)$.*

192 **Theorem 1** (Guarantees for finite classes). *Consider a triple $(\Pi, \mathcal{F}, \mathcal{Q})$ that is weakly Bellman
 193 realizable (Assumption 1); an i.i.d. dataset (Assumption 2); and the choices $\rho = c\{\log(|\mathcal{F}||\Pi||\mathcal{Q}|) +$
 194 $\log(1/\delta)\}$ and $\lambda = c'\rho/n$ for some constants c, c' . Then w.p. at least $1 - \delta$:*

- 195 • Policy evaluation: For any $\pi \in \Pi$, the estimates $(\widehat{V}_{\min}^\pi, \widehat{V}_{\max}^\pi)$ specify a confidence interval
 196 satisfying the coverage (8a) and width bounds (8b)
- 197 • Policy optimization: Any max-min policy (6) $\tilde{\pi}$ satisfies the oracle inequality (9).

198 4 Concentrability Coefficients and Test Spaces

199 In this section, we develop some connections to concentrability coefficients that have been used in
 200 past work, and discuss various choices of the test class. Like the predictor class \mathcal{Q}^π , the test class
 201 \mathcal{F}^π encodes domain knowledge, and thus its choice is delicate. Different from the predictor class,
 202 the test class does not require a ‘realizability’ condition. As a general principle, the test functions
 203 should be chosen as orthogonal as possible with respect to the Bellman residual, so as to enable
 204 rapid progress towards the solution; at the same time, they should be sufficiently “aligned” with the
 205 dataset, meaning that $\|f\|_\mu$ or its empirical counterpart $\|f\|_n$ should be large. Given a test class, each
 206 additional test function posits a new constraint which helps identify the correct predictor; at the same

²See Section B.2.1 for a precise definition of the relevant μ for a fairly general sampling model.

207 time, it increases the metric entropy (parameter ρ), which makes each individual constraints more
 208 loose. In summary, there are trade-offs to be made in the selection of the test class \mathcal{F} , much like \mathcal{Q} .
 209 In order to assess the statistical cost that we pay for off-policy data, it is natural to define the *off-policy*
 210 *cost coefficient* (OPC) as

$$K^\pi(\mathcal{C}_n^\pi, \rho, \lambda) \stackrel{def}{=} \max_{Q \in \mathcal{C}_n^\pi} \frac{|\mathbb{E}_\pi \mathcal{B}^\pi Q|^2}{(1+\lambda)\frac{\rho}{n}} = \max_{Q \in \mathcal{C}_n^\pi} \frac{\langle \mathbb{1}, \mathcal{B}^\pi Q \rangle_\mu^2}{(1+\lambda)\frac{\rho}{n}}, \quad (10)$$

211 With this notation, our off-policy width bound (8b) can be re-expressed as

$$|\widehat{V}_{\min}^\pi - \widehat{V}_{\max}^\pi| \leq 2 \frac{\sqrt{1+\lambda}}{1-\gamma} \sqrt{K^\pi \frac{\rho}{n}}, \quad (11a)$$

212 while the oracle inequality (9) for policy optimization can be re-expressed in the form

$$V^{\widehat{\pi}} \geq \max_{\pi \in \Pi} \left\{ V^\pi - \frac{\sqrt{1+\lambda}}{1-\gamma} \sqrt{K^\pi \frac{\rho}{n}} \right\}, \quad (11b)$$

213 Since $\lambda \sim \rho/n$, the factor $\sqrt{1+\lambda}$ can be bounded by a constant in the typical case $n \geq \rho$. We now
 214 offer concrete examples of the OPC, while deferring further examples to Appendix A.5.

215 4.1 Likelihood ratios

216 Our broader goal is to obtain small Bellman error along the distribution induced by π . Assume that
 217 one constructs a test function class \mathcal{F}^π of possible likelihood ratios.

218 **Proposition 1** (Likelihood ratio bounds). *Assume that for some constant b_π , the test function defined*
 219 *as $f^*(s, a) = \frac{1}{b_\pi} \frac{d_\pi(s, a)}{\mu(s, a)}$ belongs to \mathcal{F}^π and satisfies $\|f^*\|_\infty \leq 1$. Then the OPC coefficient satisfies*

$$K^\pi \leq \frac{(i) \mathbb{E}_\pi \left[\frac{d_\pi(S, A)}{\mu(S, A)} \right] + b_\pi^2 \lambda}{1+\lambda} \stackrel{(ii)}{\leq} \frac{b_\pi (1+b_\pi \lambda)}{1+\lambda} \quad (12)$$

220 The proof is in Appendix D.1. Since $\lambda = \lambda_n \rightarrow 0$ as n increases, the OPC coefficient is bounded
 221 by a multiple of the expected ratio $\mathbb{E}_\pi \left[\frac{d_\pi(S, A)}{\mu(S, A)} \right]$. Up to an additive offset, this expectation is
 222 equivalent to the χ^2 -distribution between the policy-induced occupation measure d_π and data-
 223 generating distribution μ . The concentrability coefficient can be plugged back into Eqs. (11a)
 224 and (11b) to obtain a concrete policy optimization bound. In this case, we recover a result similar to
 225 [XJ20b], but with a much milder concentrability coefficient that involves only the chosen comparator
 226 policy.

227 4.2 The error test space

228 We now turn to the discussion of a choice for the test space that extends the LSTD algorithm to
 229 non-linear spaces. A simplification to the linear setting is presented later in Section 5.

230 As is well known, the LSTD algorithm [BB96] can be seen as minimizing the Bellman error projected
 231 onto the linear prediction space \mathcal{Q} . Define the transition operator $(\mathbb{P}^\pi Q)(s, a) = \mathbb{E}_{s^+ \sim \mathbb{P}(s, a)} Q(s^+, \pi)$,
 232 and the prediction error $\epsilon = Q - Q_\star^\pi$, where Q_\star^π is a \mathcal{Q} -function from the definition of weak
 233 realizability. The Bellman error can be re-written as $\mathcal{B}^\pi Q = \mathcal{B}^\pi Q - \mathcal{B}^\pi Q_\star^\pi = (\mathcal{I} - \gamma \mathbb{P}^\pi) \epsilon$. When
 234 realizability holds, in the linear setting and at the population level, the LSTD solution seeks to satisfy
 235 the projected Bellman equations

$$\langle f, \mathcal{B}^\pi Q \rangle_\mu = 0, \quad \text{for all } f \in \mathcal{E}_\star^\pi. \quad (13)$$

236 In the linear case, \mathcal{E}_\star^π is the class of linear functions \mathcal{Q}^π used as predictors; when \mathcal{Q}^π is non-linear,
 237 we can extend the LSTD method by using the (nonlinear) error test space $\mathcal{F}^\pi = \mathcal{E}_\star^\pi = \{Q - Q_\star^\pi\}$.
 238 Since \mathcal{E}_\star^π is unknown (as it depends on the weak solution Q_\star^π), we choose instead the larger class

$$\mathcal{E}^\pi = \{Q - Q' \mid Q, Q' \in \mathcal{Q}^\pi\},$$

239 which contains \mathcal{E}_\star^π . The resulting approach can be seen as performing a projection of the Bellman
 240 operator $\mathcal{B}^\pi Q$ into the error space \mathcal{E}_\star^π , much like LSTD does in the linear setting. However, different
 241 from LSTD, our procedure returns confidence intervals as opposed to a point estimator. This choice
 242 of the test space is related to the Bubnov-Galerkin method [Rep17] for linear spaces; it selects the
 243 test space \mathcal{F}^π to be identical to the trial space \mathcal{E}_\star^π that contains all possible solution errors.

244 **Lemma 1** (OPC coefficient from prediction error). *For any test function class $\mathcal{F}^\pi \supseteq \mathcal{E}^\pi$, we have*

$$K^\pi \leq \max_{Q \in \mathcal{Q}^\pi} \left\{ \frac{\|\epsilon\|_\mu^2 + \lambda}{\|\mathbb{1}\|_\pi^2 + \lambda} \frac{\langle \mathbb{1}, \mathcal{B}^\pi Q \rangle_\mu^2}{\langle \epsilon, \mathcal{B}^\pi Q \rangle_\mu^2} \right\} = \max_{\epsilon \in \mathcal{E}_*^\pi} \left\{ \frac{\|\epsilon\|_\mu^2 + \lambda}{\|\mathbb{1}\|_\pi^2 + \lambda} \frac{\langle \mathbb{1}, (\mathcal{I} - \gamma \mathbb{P}^\pi) \epsilon \rangle_\mu^2}{\langle \epsilon, (\mathcal{I} - \gamma \mathbb{P}^\pi) \epsilon \rangle_\mu^2} \right\}. \quad (14)$$

245 The above coefficient measures the ratio between the Bellman error along the distribution of the
 246 target policy π and that projected onto the error space \mathcal{E}_*^π defined by \mathcal{Q}^π . It is a concentrability
 247 coefficient that *always* applies, as the choice of the test space does not require domain knowledge.
 248 See Appendix D.2 for the proof, and Appendix A.6 for further comments and insights, as well as a
 249 simplification in the special case of Bellman closure.

250 4.3 The Bellman test space

251 In the prior section we controlled the projected Bellman error. Another longstanding approach in
 252 reinforcement learning is to control the Bellman error itself, for example by minimizing the squared
 253 Bellman residual. In general, this cannot be done if only an offline dataset is available due to the well
 254 known *double sampling* issue. However, in some cases we can use an helper class to try to capture
 255 the Bellman error. Such class needs to be a superset of the class of *Bellman test functions* given by

$$\mathcal{F}_\pi^{\mathcal{B}} \stackrel{\text{def}}{=} \{\mathcal{B}^\pi Q \mid Q \in \mathcal{Q}^\pi\}. \quad (15)$$

256 Any test class that contains the above allows us to control the Bellman residual, as we show next.

257 **Lemma 2** (Bellman Test Functions). *For any test function class \mathcal{F}^π that contains $\mathcal{F}_\pi^{\mathcal{B}}$, we have*

$$\|\mathcal{B}^\pi Q\|_\mu \leq c_1 \sqrt{\frac{\rho}{n}} \quad \text{for any } Q \in \mathcal{C}_n^\pi(\mathcal{F}^\pi). \quad (16a)$$

258 *Moreover, the off-policy cost coefficient is upper bounded as*

$$K^\pi \stackrel{(i)}{\leq} c_1 \sup_{Q \in \mathcal{Q}^\pi} \frac{\langle \mathbb{1}, \mathcal{B}^\pi Q \rangle_\mu^2}{\|\mathcal{B}^\pi Q\|_\mu^2} \stackrel{(ii)}{\leq} c_1 \sup_{Q \in \mathcal{Q}^\pi} \frac{\|\mathcal{B}^\pi Q\|_\pi^2}{\|\mathcal{B}^\pi Q\|_\mu^2} \stackrel{(iii)}{\leq} c_1 \sup_{(s,a)} \frac{d_\pi(s,a)}{\mu(s,a)}. \quad (16b)$$

259 See Appendix D.4 for the proof of this claim.

260 Consequently, whenever the test class includes the Bellman test functions, the off-policy cost
 261 coefficient is at most the ratio between the squared Bellman residuals along the data generating
 262 distribution and the target distribution. If Bellman closure holds, then the prediction error space
 263 \mathcal{E}^π introduced in Section 4.2 contains the Bellman test functions: for $Q \in \mathcal{Q}^\pi$, we can write
 264 $\mathcal{B}^\pi Q = Q - \mathcal{T}^\pi Q \in \mathcal{E}^\pi$. This fact allows us to recover a result in the recent paper [XCJ⁺21] in the
 265 special case of Bellman closure, although the approach presented here is more general.

266 4.4 Combining test spaces

267 Often, it is natural to construct a test space that is a union of several simpler classes. A simple but
 268 valuable observation is that the resulting procedure inherits the best of the OPC coefficients. Suppose
 269 that we are given a collection $\{\mathcal{F}_m^\pi\}_{m=1}^M$ of M different test function classes, and define the union
 270 $\mathcal{F}^\pi = \bigcup_{m=1}^M \mathcal{F}_m^\pi$. For each $m = 1, \dots, M$, let K_m^π be the OPC coefficient defined by the function
 271 class \mathcal{F}_m^π and radius ρ , and let $K^\pi(\mathcal{F})$ be the OPC coefficient associated with the full class. Then we
 272 have the following guarantee:

273 **Lemma 3** (Multiple test classes). $K^\pi(\mathcal{F}) \leq \min_{m=1, \dots, M} K_m^\pi$.

274 This guarantee is a straightforward consequence of our construction of the feasibility sets: in particular,
 275 we have $\mathcal{C}_n^\pi(\mathcal{F}) = \bigcap_{m=1}^M \mathcal{C}_n^\pi(\mathcal{F}_m)$, and consequently, by the variational definition of the off-policy
 276 cost coefficient $K^\pi(\mathcal{F})$ as optimization over $\mathcal{C}_n^\pi(\mathcal{F})$, the bound (3) follows. In words, when multiple
 277 test spaces are combined, then our algorithms inherit the best (smallest) OPC coefficient over all
 278 individual test spaces. While this behavior is attractive, one must note that there is a statistical cost to
 279 using a union of test spaces: the choice of ρ scales as a function of \mathcal{F} via its metric entropy. This
 280 increase in ρ must be balanced with the benefits of using multiple test spaces.³

³For space reasons, we defer to Appendix A.7 an application in which we construct a test function space as a union of subclasses, and thereby obtain a method that automatically leverages Bellman closure when it holds, falls back to importance sampling if closure fails, and falls back to a worst-case bound in general.

281 5 Linear Setting

282 In this section, we turn to a detailed analysis of our estimators using function classes that are linear
 283 in a feature map. Let $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ be a given feature map, and consider linear expansions

284 $g_w(s, a) \stackrel{def}{=} \langle w, \phi(s, a) \rangle = \sum_{j=1}^d w_j \phi_j(s, a)$. The class of *linear functions* takes the form

$$\mathcal{L} \stackrel{def}{=} \{(s, a) \mapsto g_w(s, a) \mid w \in \mathbb{R}^d, \|w\|_2 \leq 1\}. \quad (17)$$

285 Throughout our analysis, we assume that $\|\phi(s, a)\|_2 \leq 1$ for all state-action pairs.

286 Following the approach in Section 4.2, which is based on the LSTD method, we should choose the
 287 test function class $\mathcal{F}^\pi = \mathcal{L}$, as in the linear case the prediction error is linear.

288 In order to obtain a computationally efficient implementation, we need to use a test class that is a
 289 “simpler” subset of \mathcal{L} . In particular, for linear functions, it is not hard to show that the estimates
 290 \widehat{V}_{\min}^π and \widehat{V}_{\max}^π from equation (5) can be computed by solving a quadratic program, with two linear
 291 constraints for each test function. (See Appendix A.8 for the details.) Consequently, the computational
 292 complexity scales linearly with the number of test functions. Thus, if we restrict ourselves to a finite
 293 test class contained within \mathcal{L} , we will obtain a computationally efficient approach.

294 5.1 A computationally friendly test class and OPC coefficients

295 Define the empirical covariance matrix $\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \phi_i \phi_i^T$ where $\phi_i \stackrel{def}{=} \phi(s_i, a_i)$. Let $\{\widehat{u}_j\}_{j=1}^d$ be
 296 the eigenvectors of empirical covariance matrix $\widehat{\Sigma}$, and suppose that they are normalized to have unit
 297 ℓ_2 -norm. We use these normalized eigenvectors to define the finite test class

$$\widetilde{\mathcal{F}}^\pi \stackrel{def}{=} \{f_j, j = 1, \dots, d\} \quad \text{where } f_j(s, a) \stackrel{def}{=} \langle \widehat{u}_j, \phi(s, a) \rangle \quad (18)$$

298 A few observations are in order:

- 299 • This test class has only d functions, so that our QP implementation has $2d$ constraints, and can
 300 be solved in polynomial time. (Again, see Appendix A.8 for details.)
- 301 • Since $\widetilde{\mathcal{F}}^\pi$ is a subset of \mathcal{L} the choice of radius $\rho = c(\frac{d}{n} + \log 1/\delta)$ is valid for some constant c .

302 **Concentrability:** When weak Bellman closure does not hold, then our analysis needs to take into
 303 account how errors propagate via the dynamics. In particular, we define the *next-state feature*
 304 *extractor* $\phi^{+\pi}(s, a) \stackrel{def}{=} \mathbb{E}_{s^+ \sim \mathbb{P}(s, a)} \phi(s^+, \pi)$, along with the population covariance matrix $\Sigma \stackrel{def}{=} \mathbb{E}_\mu[\phi(s, a) \phi^\top(s, a)]$, and its λ -regularized version $\Sigma_\lambda \stackrel{def}{=} \Sigma + \lambda I$. We also define the matrices

$$\Sigma^{+\pi} \stackrel{def}{=} \mathbb{E}_\mu[\phi(\phi^{+\pi})^\top], \quad \Sigma_{\lambda, \text{Boot}}^{+\pi} \stackrel{def}{=} (\Sigma_\lambda^{\frac{1}{2}} - \gamma \Sigma_\lambda^{-\frac{1}{2}} \Sigma^{+\pi})^\top (\Sigma_\lambda^{\frac{1}{2}} - \gamma \Sigma_\lambda^{-\frac{1}{2}} \Sigma^{+\pi}).$$

306 The matrix $\Sigma^{+\pi}$ is the cross-covariance between successive states, whereas the matrix $\Sigma_{\lambda, \text{Boot}}^{+\pi}$ is a
 307 suitably renormalized and symmetrized version of the matrix $\Sigma_\lambda^{\frac{1}{2}} - \gamma \Sigma_\lambda^{-\frac{1}{2}} \Sigma^{+\pi}$, which arises naturally
 308 from the policy evaluation equation. We refer to quantities that contain evaluations at the next-state
 309 (e.g., $\phi^{+\pi}$) as bootstrapping terms, and now bound the OPC coefficient in the presence of such terms:
 310 **Proposition 2** (OPC bounds with bootstrapping). *Under weak realizability, we have*

$$K^\pi(\widetilde{\mathcal{F}}^\pi) \leq c d \|\mathbb{E}_\pi[\phi - \gamma \phi^{+\pi}]\|_{(\Sigma_{\lambda, \text{Boot}}^{+\pi})^{-1}}^2 \quad \text{with probability at least } 1 - \delta. \quad (19)$$

311 See Appendix E.1 for the proof. The bound (19) takes a familiar form, as it involves the same
 312 matrices used to define the LSTD solution. This is expected, as our approach here is essentially
 313 equivalent to the LSTD method; the difference is that LSTD only gives a point estimate as opposed to
 314 the confidence intervals that we present here; however, they are both derived from the same principle,
 315 namely from the Bellman equations projected along the predictor (error) space.

316 The bound quantifies how the feature extractor ϕ together with the bootstrapping term $\phi^{+\pi}$, averaged
 317 along the target policy π , interact with the covariance matrix with bootstrapping $\Sigma_{\lambda, \text{Boot}}^{+\pi}$. It is an
 318 approximation to the OPC coefficient bound derived in Lemma 1. The bootstrapping terms capture
 319 the temporal difference correlations that can arise in reinforcement learning when strong assumptions
 320 like Bellman closure do not hold. As a consequence, such an OPC coefficient being small is a
 321 *sufficient* condition for reliable off-policy prediction. This bound on the OPC coefficient always
 322 applies, and it reduces to the simpler one (20) when weak Bellman closure holds, with no need to
 323 inform the algorithm of the simplified setting; see Appendix E.3 for the proof.

324 **Proposition 3** (OPC bounds under weak Bellman Closure). *Under Bellman closure, we have*

$$K^\pi(\tilde{\mathcal{F}}^\pi) \leq c d \|\mathbb{E}_\pi \phi\|_{\Sigma_\lambda^{-1}}^2 \quad \text{with probability at least } 1 - \delta. \quad (20)$$

325 5.2 Actor-critic scheme for policy optimization

326 Having described a practical procedure to compute \widehat{V}_{\min}^π , we now turn to the computation of the
327 max-min estimator for policy optimization. We define the *soft-max policy class*

$$\Pi_{\text{lin}} \stackrel{\text{def}}{=} \left\{ (s, a) \mapsto \frac{e^{\langle \phi(s, a), \theta \rangle}}{\sum_{a' \in \mathcal{A}} e^{\langle \phi(s, a'), \theta \rangle}} \mid \|\theta\|_2 \leq T, \theta \in \mathbb{R}^d \right\}. \quad (21)$$

328 In order to compute the max-min solution (6) over this policy class, we implement an actor-critic
329 method, in which the actor performs a variant of mirror descent.⁴

330 • At each iteration $t = 1, \dots, T$, the policy $\pi_t \in \Pi_{\text{lin}}$ can be identified with a parameter $\theta_t \in \mathbb{R}^d$.
331 The sequence is initialized with $\theta_1 = 0$.

332 • Using the finite test function class (18) based on normalized eigenvectors, the pessimistic
333 value estimate $\widehat{V}_{\min}^{\pi_t}$ is computed by solving a quadratic program, as previously described. This
334 computation returns the weight vector w_t of the associated optimal action-value function.

335 • Using the action-value vector w_t , we update the actor's parameter as

$$\theta_{t+1} = \theta_t + \eta w_t \quad \text{where } \eta = \sqrt{\frac{\log|\mathcal{A}|}{2T}} \text{ is a stepsize parameter.} \quad (22)$$

336 We now state a guarantee on the behavior of this procedure, based on two OPC coefficients:

$$K_{(1)}^{\tilde{\pi}} = d \|\mathbb{E}_{\tilde{\pi}} \phi\|_{\Sigma_\lambda^{-1}}^2, \quad \text{and} \quad K_{(2)}^{\tilde{\pi}} = d \sup_{\pi \in \Pi} \left\{ \|\mathbb{E}_{\tilde{\pi}}[\phi - \gamma \phi^{+\pi}]\|_{(\Sigma_\lambda^{+\pi})^{-1}}^2 \right\}. \quad (23)$$

337 Moreover, in making the following assertion, we assume that every weak solution Q_\star^π can be evaluated
338 against the distribution of a comparator policy $\tilde{\pi} \in \Pi$, i.e., $\langle \mathbf{1}, \mathcal{B}^\pi Q_\star^\pi \rangle_{\tilde{\pi}} = 0$ for all $\pi \in \Pi$. (This
339 assumption is still weaker than strong realizability).

340 **Theorem 2** (Approximate Guarantees for Linear Soft-Max Optimization). *Under the above con-*
341 *ditions, running the procedure for T rounds returns a policy sequence $\{\pi_t\}_{t=1}^T$ such that, for any*
342 *comparator policy $\tilde{\pi} \in \Pi$,*

$$\frac{1}{T} \sum_{t=1}^T \{V^{\tilde{\pi}} - V^{\pi_t}\} \leq \frac{c_1}{1-\gamma} \left\{ \underbrace{\sqrt{\frac{\log|\mathcal{A}|}{T}}}_{\text{Optimization error}} + \underbrace{\sqrt{K_{(\cdot)}^{\tilde{\pi}} \frac{d \log(nT) + \log\left(\frac{n}{\delta}\right)}{n}}}_{\text{Statistical error}} \right\}, \quad (24)$$

343 *with probability at least $1 - \delta$. This bound always holds with $K_{(\cdot)}^{\tilde{\pi}} = K_{(2)}^{\tilde{\pi}}$, and moreover, it holds*
344 *with $K_{(\cdot)}^{\tilde{\pi}} = K_{(1)}^{\tilde{\pi}}$ when weak Bellman closure is in force.*

345 See Appendix F for the proof. Whenever Bellman closure holds, the result automatically inherits the
346 more favorable concentrability coefficient $K_{(2)}^{\tilde{\pi}}$, as originally derived in Proposition 3. The resulting
347 bound is only \sqrt{d} worse than the lower bound recently established in the paper [ZWB21]. However,
348 the method proposed here is robust, in that it provides guarantees even when Bellman closure does not
349 hold. In this case, we have a guarantee in terms of the OPC coefficient $K_{(1)}^{\tilde{\pi}}$. Note that it is a uniform
350 version of the one derived previously in Proposition 2, in that there is an additional supremum over
351 the policy class. This supremum arises due to the use of gradient-based method, which implicitly
352 searches over policies in bootstrapping terms; see Appendix A.9 for a more detailed discussion of
353 this issue.

354 References

355 [AMS07] András Antos, Rémi Munos, and Csaba Szepesvári. Fitted Q-iteration in continuous
356 action-space MDPs. 2007.

⁴Strictly speaking, it is mirror ascent, but we use the conventional terminology.

- 357 [ASM08] András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies
358 with Bellman-residual minimization based fitted policy iteration and a single sample
359 path. *Machine Learning*, 71(1):89–129, 2008.
- 360 [ASN20] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic per-
361 spective on offline reinforcement learning. In *International Conference on Machine*
362 *Learning*, pages 104–114. PMLR, 2020.
- 363 [BB96] Steven J Bradtke and Andrew G Barto. Linear least-squares algorithms for temporal
364 difference learning. *Machine learning*, 22(1-3):33–57, 1996.
- 365 [Ber95a] D. P. Bertsekas. *Dynamic programming and stochastic control*, volume 1. Athena
366 Scientific, Belmont, MA, 1995.
- 367 [Ber95b] D.P. Bertsekas. *Dynamic programming and stochastic control*, volume 2. Athena
368 Scientific, Belmont, MA, 1995.
- 369 [BGB20] Jacob Buckman, Carles Gelada, and Marc G Bellemare. The importance of pessimism
370 in fixed-dataset policy optimization. *arXiv preprint arXiv:2009.06799*, 2020.
- 371 [BLL⁺11] Alina Beygelzimer, John Langford, Lihong Li, Lev Reyzin, and Robert Schapire.
372 Contextual bandit algorithms with supervised learning guarantees. In *Proceedings of*
373 *the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages
374 19–26. JMLR Workshop and Conference Proceedings, 2011.
- 375 [BT96] Dimitri P Bertsekas and John N Tsitsiklis. *Neuro-dynamic programming*. Athena
376 Scientific, 1996.
- 377 [CJ19] Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement
378 learning. In *International Conference on Machine Learning*, pages 1042–1051, 2019.
- 379 [CQ22] Xiaohong Chen and Zhengling Qi. On well-posedness and minimax optimal rates
380 of nonparametric q-function estimation in off-policy evaluation. *arXiv preprint*
381 *arXiv:2201.06169*, 2022.
- 382 [DJL21] Yaqi Duan, Chi Jin, and Zhiyuan Li. Risk bounds and rademacher complexity in batch
383 reinforcement learning. *arXiv preprint arXiv:2103.13883*, 2021.
- 384 [DW20] Yaqi Duan and Mengdi Wang. Minimax-optimal off-policy evaluation with linear
385 function approximation. *arXiv preprint arXiv:2002.09516*, 2020.
- 386 [Eva10] Lawrence C Evans. *Partial differential equations*, volume 19. American Mathematical
387 Soc., 2010.
- 388 [FCG18] Mehrdad Farajtabar, Yinlam Chow, and Mohammad Ghavamzadeh. More robust doubly
389 robust off-policy evaluation. In *International Conference on Machine Learning*, pages
390 1447–1456. PMLR, 2018.
- 391 [FGSM16] Amir-massoud Farahmand, Mohammad Ghavamzadeh, Csaba Szepesvári, and Shie
392 Mannor. Regularized policy iteration with nonparametric function spaces. *The Journal*
393 *of Machine Learning Research*, 17(1):4809–4874, 2016.
- 394 [FKSLX21] Dylan J. Foster, Akshay Krishnamurthy, David Simchi-Levi, and Yunzong Xu. Offline
395 reinforcement learning: Fundamental barriers for value function approximation, 2021.
- 396 [Fle84] Clive AJ Fletcher. Computational Galerkin methods. In *Computational Galerkin*
397 *methods*, pages 72–85. Springer, 1984.
- 398 [FRTL20] Yihao Feng, Tongzheng Ren, Ziyang Tang, and Qiang Liu. Accountable off-policy
399 evaluation with kernel bellman statistics. In *International Conference on Machine*
400 *Learning*, pages 3102–3111. PMLR, 2020.
- 401 [FSM10] Amir-massoud Farahmand, Csaba Szepesvári, and Rémi Munos. Error propagation for
402 approximate policy and value iteration. In *Advances in Neural Information Processing*
403 *Systems (NIPS)*, 2010.

- 404 [Gal15] Boris Grigoryevich Galerkin. Series solution of some problems of elastic equilibrium
405 of rods and plates. *Vestnik inzhenerov i tekhnikov*, 19(7):897–908, 1915.
- 406 [Haz21] Elad Hazan. Introduction to online convex optimization, 2021.
- 407 [HJD⁺21] Botao Hao, Xiang Ji, Yaqi Duan, Hao Lu, Csaba Szepesvári, and Mengdi
408 Wang. Bootstrapping statistical inference for off-policy evaluation. *arXiv preprint*
409 *arXiv:2102.03607*, 2021.
- 410 [JGS⁺19] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata
411 Lapedriza, Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch
412 deep reinforcement learning of implicit human preferences in dialog. *arXiv preprint*
413 *arXiv:1907.00456*, 2019.
- 414 [JH20] Nan Jiang and Jiawei Huang. Minimax value interval for off-policy evaluation and
415 policy optimization. *arXiv preprint arXiv:2002.02081*, 2020.
- 416 [JKA⁺17] Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E.
417 Schapire. Contextual decision processes with low Bellman rank are PAC-learnable.
418 In Doina Precup and Yee Whye Teh, editors, *International Conference on Machine*
419 *Learning (ICML)*, volume 70 of *Proceedings of Machine Learning Research*, pages
420 1704–1713, International Convention Centre, Sydney, Australia, 06–11 Aug 2017.
421 PMLR.
- 422 [JL16] Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement
423 learning. In *International Conference on Machine Learning*, pages 652–661. PMLR,
424 2016.
- 425 [JYW21] Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline
426 rl? In *International Conference on Machine Learning*, pages 5084–5096. PMLR, 2021.
- 427 [K⁺03] Sham Machandranath Kakade et al. *On the sample complexity of reinforcement learning*.
428 PhD thesis, University of London London, England, 2003.
- 429 [KFTL19] Aviral Kumar, Justin Fu, George Tucker, and Sergey Levine. Stabilizing off-policy
430 q-learning via bootstrapping error reduction. *arXiv preprint arXiv:1906.00949*, 2019.
- 431 [KHSL21] Aviral Kumar, Joey Hong, Anikait Singh, and Sergey Levine. Should i run offline
432 reinforcement learning or behavioral cloning? In *Deep RL Workshop NeurIPS 2021*,
433 2021.
- 434 [KRNJ20] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims.
435 Morel: Model-based offline reinforcement learning. *arXiv preprint arXiv:2005.05951*,
436 2020.
- 437 [KU19] Nathan Kallus and Masatoshi Uehara. Efficiently breaking the curse of horizon in off-
438 policy evaluation with double reinforcement learning. *arXiv preprint arXiv:1909.05850*,
439 2019.
- 440 [LLTZ18] Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon:
441 Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing*
442 *Systems*, pages 5356–5366, 2018.
- 443 [LP03] Michail G Lagoudakis and Ronald Parr. Least-squares policy iteration. *Journal of*
444 *machine learning research*, 4(Dec):1107–1149, 2003.
- 445 [LSAB20] Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Provably
446 good batch reinforcement learning without great exploration. *arXiv preprint*
447 *arXiv:2007.08202*, 2020.
- 448 [LTDC19] Romain Laroche, Paul Trichelair, and Remi Tachet Des Combes. Safe policy improve-
449 ment with baseline bootstrapping. In *International Conference on Machine Learning*,
450 pages 3652–3661. PMLR, 2019.

- 451 [LTND21] Jonathan N Lee, George Tucker, Ofir Nachum, and Bo Dai. Model selection in batch
452 policy optimization. *arXiv preprint arXiv:2112.12320*, 2021.
- 453 [MS08] Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal*
454 *of Machine Learning Research*, 9(May):815–857, 2008.
- 455 [Mun03] Rémi Munos. Error bounds for approximate policy iteration. In *ICML*, volume 3, pages
456 560–567, 2003.
- 457 [Mun05] Rémi Munos. Error bounds for approximate value iteration. In *AAAI Conference on*
458 *Artificial Intelligence (AAAI)*, 2005.
- 459 [ND20] Ofir Nachum and Bo Dai. Reinforcement learning via Fenchel-Rockafellar duality.
460 *arXiv preprint arXiv:2001.01866*, 2020.
- 461 [NDGL20] Ashvin Nair, Murtaza Dalal, Abhishek Gupta, and Sergey Levine. Accelerating online
462 reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- 463 [NDK⁺19] Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans.
464 Algaedice: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*,
465 2019.
- 466 [Pre00] Doina Precup. Eligibility traces for off-policy policy evaluation. *Computer Science*
467 *Department Faculty Publication Series*, page 80, 2000.
- 468 [Put94] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Pro-*
469 *gramming*. John Wiley & Sons, Inc., New York, NY, USA, 1994.
- 470 [Rep17] Sergey Repin. One hundred years of the Galerkin method. *Computational Methods in*
471 *Applied Mathematics*, 17(3):351–357, 2017.
- 472 [RZM⁺21] Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging
473 offline reinforcement learning and imitation learning: A tale of pessimism. *arXiv*
474 *preprint arXiv:2103.12021*, 2021.
- 475 [SB18] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT
476 Press, 2018.
- 477 [SSB⁺20] Noah Y Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki,
478 Michael Neunert, Thomas Lampe, Roland Hafner, Nicolas Heess, and Martin Riedmiller.
479 Keep doing what worked: Behavioral modelling priors for offline reinforcement learning.
480 *arXiv preprint arXiv:2002.08396*, 2020.
- 481 [TB16] Philip Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for
482 reinforcement learning. In *International Conference on Machine Learning*, pages
483 2139–2148, 2016.
- 484 [TFL⁺19] Ziyang Tang, Yihao Feng, Lihong Li, Dengyong Zhou, and Qiang Liu. Doubly robust
485 bias reduction in infinite horizon off-policy estimation. *arXiv preprint arXiv:1910.07186*,
486 2019.
- 487 [UHJ20] Masatoshi Uehara, Jiawei Huang, and Nan Jiang. Minimax weight and q-function
488 learning for off-policy evaluation. In *International Conference on Machine Learning*,
489 pages 9659–9668. PMLR, 2020.
- 490 [UIJ⁺21] Masatoshi Uehara, Masaaki Imaizumi, Nan Jiang, Nathan Kallus, Wen Sun, and
491 Tengyang Xie. Finite sample analysis of minimax offline reinforcement learning:
492 Completeness, fast rates and first-order efficiency. *arXiv preprint arXiv:2102.02981*,
493 2021.
- 494 [US21] Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning
495 under partial coverage, 2021.

- 496 [VJY21] Cameron Voloshin, Nan Jiang, and Yisong Yue. Minimax model learning. In *International Conference on Artificial Intelligence and Statistics*, pages 1612–1620. PMLR, 497 2021. 498
- 499 [Wai19] Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, vol- 500 ume 48. Cambridge University Press, 2019.
- 501 [WFK20] Ruosong Wang, Dean P Foster, and Sham M Kakade. What are the statistical limits of 502 offline rl with linear function approximation? *arXiv preprint arXiv:2010.11895*, 2020.
- 503 [WNŻ⁺20] Ziyu Wang, Alexander Novikov, Konrad Żołna, Jost Tobias Springenberg, Scott Reed, 504 Bobak Shahriari, Noah Siegel, Josh Merel, Caglar Gulcehre, Nicolas Heess, et al. Critic 505 regularized regression. *arXiv preprint arXiv:2006.15134*, 2020.
- 506 [WTN19] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement 507 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 508 [XCJ⁺21] Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. 509 Bellman-consistent pessimism for offline reinforcement learning. *arXiv preprint 510 arXiv:2106.06926*, 2021.
- 511 [XJ20a] Tengyang Xie and Nan Jiang. Batch value-function approximation with only realizabil- 512 ity. *arXiv preprint arXiv:2008.04990*, 2020.
- 513 [XJ20b] Tengyang Xie and Nan Jiang. Q* approximation schemes for batch reinforcement 514 learning: A theoretical comparison. volume 124 of *Proceedings of Machine Learning 515 Research*, pages 550–559, Virtual, 03–06 Aug 2020. PMLR.
- 516 [XMW19] Tengyang Xie, Yifei Ma, and Yu-Xiang Wang. Towards optimal off-policy evaluation 517 for reinforcement learning with marginalized importance sampling. In *Advances in 518 Neural Information Processing Systems*, pages 9668–9678, 2019.
- 519 [YB10] H. Yu and D. P. Bertsekas. Error bounds for approximations from projected linear 520 equations. *Mathematics of Operations Research*, 35(2):306–329, 2010.
- 521 [YBW20] Ming Yin, Yu Bai, and Yu-Xiang Wang. Near optimal provable uniform convergence 522 in off-policy evaluation for reinforcement learning. *arXiv preprint arXiv:2007.03760*, 523 2020.
- 524 [YND⁺20] Mengjiao Yang, Ofir Nachum, Bo Dai, Lihong Li, and Dale Schuurmans. Off-policy 525 evaluation via the regularized lagrangian. *arXiv preprint arXiv:2007.03438*, 2020.
- 526 [YQCC21] Chao-Han Huck Yang, Zhengling Qi, Yifan Cui, and Pin-Yu Chen. Pessimistic model 527 selection for offline deep reinforcement learning, 2021.
- 528 [YTY⁺20] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, 529 Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. *arXiv 530 preprint arXiv:2005.13239*, 2020.
- 531 [YW20] Ming Yin and Yu-Xiang Wang. Asymptotically efficient off-policy evaluation for 532 tabular reinforcement learning. In *International Conference on Artificial Intelligence 533 and Statistics*, pages 3948–3958. PMLR, 2020.
- 534 [YW21] Ming Yin and Yu-Xiang Wang. Towards instance-optimal offline reinforcement learning 535 with pessimism. *arXiv preprint arXiv:2110.08695*, 2021.
- 536 [YWDW] Ming Yin, Yu-Xiang Wang, Yaqi Duan, and Mengdi Wang. Near-optimal offline 537 reinforcement learning with linear representation: Leveraging variance information 538 with pessimism.
- 539 [Zan20] Andrea Zanette. Exponential lower bounds for batch reinforcement learning: Batch rl 540 can be exponentially harder than online RL. *arXiv preprint arXiv:2012.08005*, 2020.
- 541 [ZDLS20] Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. Gendice: Generalized offline 542 estimation of stationary values. *arXiv preprint arXiv:2002.09072*, 2020.

- 543 [ZHH⁺22] Wenhao Zhan, Baihe Huang, Audrey Huang, Nan Jiang, and Jason D Lee. Offline rein-
544 forcement learning with realizability and single-policy concentrability. *arXiv preprint*
545 *arXiv:2202.04634*, 2022.
- 546 [ZJZ21] Zihan Zhang, Xiangyang Ji, and Yuan Zhou. Almost optimal batch-regret tradeoff for
547 batch linear contextual bandits, 2021.
- 548 [ZLW20] Shangdong Zhang, Bo Liu, and Shimon Whiteson. Gradientdice: Rethinking general-
549 ized offline estimation of stationary values. In *International Conference on Machine*
550 *Learning*, pages 11194–11203. PMLR, 2020.
- 551 [ZSU⁺22] Xuezhou Zhang, Yuda Song, Masatoshi Uehara, Mengdi Wang, Alekh Agarwal, and
552 Wen Sun. Efficient reinforcement learning in block MDPs: A model-free representation
553 learning approach, 2022.
- 554 [ZWB21] Andrea Zanette, Martin J Wainwright, and Emma Brunskill. Provable benefits of actor-
555 critic methods for offline reinforcement learning. *arXiv preprint arXiv:2108.08812*,
556 2021.

557 **Checklist**

- 558 1. For all authors...
- 559 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
560 contributions and scope? [Yes]
- 561 (b) Did you describe the limitations of your work? [Yes] Limitations include the choice of
562 the test functions
- 563 (c) Did you discuss any potential negative societal impacts of your work? [N/A] The work
564 is theoretical in nature, no negative societal impacts anticipated
- 565 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
566 them? [Yes]
- 567 2. If you are including theoretical results...
- 568 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 569 (b) Did you include complete proofs of all theoretical results? [Yes]
- 570 3. If you ran experiments...
- 571 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
572 mental results (either in the supplemental material or as a URL)? [N/A]
- 573 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
574 were chosen)? [N/A]
- 575 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
576 ments multiple times)? [N/A]
- 577 (d) Did you include the total amount of compute and the type of resources used (e.g., type
578 of GPUs, internal cluster, or cloud provider)? [N/A]
- 579 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 580 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 581 (b) Did you mention the license of the assets? [N/A]
- 582 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
583
- 584 (d) Did you discuss whether and how consent was obtained from people whose data you're
585 using/curating? [N/A]
- 586 (e) Did you discuss whether the data you are using/curating contains personally identifiable
587 information or offensive content? [N/A]
- 588 5. If you used crowdsourcing or conducted research with human subjects...
- 589 (a) Did you include the full text of instructions given to participants and screenshots, if
590 applicable? [N/A]
- 591 (b) Did you describe any potential participant risks, with links to Institutional Review
592 Board (IRB) approvals, if applicable? [N/A]
- 593 (c) Did you include the estimated hourly wage paid to participants and the total amount
594 spent on participant compensation? [N/A]

595	Contents	
596	1 Introduction	1
597	2 Background and set-up	2
598	2.1 Markov decision processes and Bellman errors	3
599	2.2 Function Spaces and Weak Representation	3
600	3 Policy Estimates via the Weak Bellman Equations	4
601	3.1 Weak Bellman equations, empirical approximations and confidence intervals	4
602	3.2 High-probability guarantees	5
603	4 Concentrability Coefficients and Test Spaces	5
604	4.1 Likelihood ratios	6
605	4.2 The error test space	6
606	4.3 The Bellman test space	7
607	4.4 Combining test spaces	7
608	5 Linear Setting	8
609	5.1 A computationally friendly test class and OPC coefficients	8
610	5.2 Actor-critic scheme for policy optimization	9
611	A Additional Discussion and Results	18
612	A.1 Bellman Residual Orthogonalization	18
613	A.2 Comparison with Weight Learning Methods	18
614	A.3 Additional Literature	19
615	A.4 Definition of Weak Bellman Closure	19
616	A.5 Additional results on the concentrability coefficients	20
617	A.5.1 Testing with the identity function	20
618	A.5.2 Mixture distributions	20
619	A.5.3 Bellman Rank for off-policy evaluation	21
620	A.6 Further comments on the prediction error test space	22
621	A.7 From Importance Sampling to Bellman Closure	22
622	A.8 Implementation for Off-Policy Predictions	23
623	A.9 Discussion of Linear Approximate Optimization	24
624	B General Guarantees	25
625	B.1 A deterministic guarantee	25
626	B.2 Some high-probability guarantees	25
627	B.2.1 A model for data generation	26
628	B.2.2 A general guarantee	26
629	B.2.3 Some corollaries	28

630	C Main Proofs	30
631	C.1 Proof of Proposition 4	30
632	C.1.1 Proof of policy evaluation claims	30
633	C.1.2 Proof of policy optimization claims	30
634	C.2 Proof of Lemma 9	31
635	C.3 Proof of Theorem 3	31
636	C.4 Proof of the sandwich relation (51a)	32
637	C.5 Proof of the uniform upper bound (51b)	33
638	C.6 Proofs of supporting lemmas	34
639	C.6.1 Proof of Lemma 11	34
640	C.6.2 Proof of Lemma 12	35
641	C.6.3 Proof of Lemma 10	36
642	D Proofs for Section 4 and Appendix A.5	37
643	D.1 Proof of Proposition 1	37
644	D.2 Proof of Lemma 1	37
645	D.3 Proof of Lemma 7	37
646	D.4 Proof of Lemma 2	38
647	E Proofs for the Linear Setting	39
648	E.1 Proof of Proposition 2	39
649	E.2 Proof of Lemma 15	39
650	E.3 Proof of Proposition 3	40
651	E.4 Proof of Appendix E.4	41
652	F Proof of Theorem 2	42
653	F.1 Adversarial MDPs	42
654	F.2 Equivalence of Updates	42
655	F.3 Mirror Descent on Adversarial MDPs	43
656	F.4 Pessimism: Bound on $\mathbb{E}_{\pi_t} R_t$	43
657	F.5 Concentrability: Bound on $\mathbb{E}_{\tilde{\pi}} R_t$	44
658	F.6 Proof of Lemma 19	44

659 **A Additional Discussion and Results**

660 **A.1 Bellman Residual Orthogonalization**

661 Suppose that our goal is to estimate the action-value function Q^π of a given policy π . This function is
 662 known to be a fixed point of the Bellman evaluation operator \mathcal{T}^π associated with the policy π . Thus,
 663 when the MDP is known, one option is to (approximately) solve the Bellman evaluation equations
 664 $Q(s, a) = (\mathcal{T}^\pi Q)(s, a)$ for all state-action pairs. However, even if function approximation for Q is
 665 implemented, it is still difficult to directly solve these equations if the state-action space is sufficiently
 666 complex.

667 This observation motivates the strategy taken in this paper: instead of enforcing the Bellman equations
 668 for all state-action pairs, suppose that we do so only in an average sense, and with respect to a certain
 669 set of functions. More formally, a *test function* is a mapping from the state-action space to the real
 670 line; any such function serves to enforce the Bellman equations in an average sense in the following
 671 way. Let \mathcal{F}^π denote some user-prescribed class of test functions, which we refer to as the *test space*.
 672 Then for a given measure μ , we require only that the action-value function Q^π satisfy the integral
 673 constraints

$$\langle f, Q - \mathcal{T}^\pi(Q) \rangle_\mu \stackrel{\text{def}}{=} \int f(s, a)[Q(s, a) - (\mathcal{T}^\pi Q)(s, a)]d\mu = 0, \quad \text{for all } f \in \mathcal{F}^\pi. \quad (25)$$

674 We refer to this design principle as *Bellman residual orthogonalization*, because it requires the
 675 Bellman error function to be orthogonal to a set of test functions, as measured under the $L^2(\mu)$ inner
 676 product. Of course, by enlarging the test space \mathcal{F}^π , the Bellman error is required to be orthogonal to
 677 more test functions, and it will ultimately be zero if enough test functions are added as constraints.
 678 But at the same time, as shown by our analysis, any such enlargement has both computational and
 679 statistical costs, so there are tradeoffs to be understood.

680 In numerical analysis, especially in solving partial differential equations, the design principle (25) is
 681 called the weak or variational formulation (e.g., [Eva10]), and its solutions are referred to as weak
 682 solutions. Here we are advocating a *weak formulation* of the Bellman equations. Of course, the
 683 constraints (25) are necessary but not sufficient: the *weak (Bellman) solutions* need not solve the
 684 Bellman equations. However, whenever we need to learn based on a limited dataset, it is unreasonable
 685 to satisfy the Bellman equations everywhere; instead, by choosing the test space appropriately, we can
 686 seek to satisfy the Bellman equations over regions of the state-action space that are most important.
 687 In some cases, the formulation (25) can be fruitfully viewed as a type of Galerkin approximation
 688 (e.g., [Gal15, Fle84]) to the Bellman equations. For example, when both the test functions and
 689 Q -value functions belong to some linear space (and the empirical constraints are enforced exactly),
 690 then the weak formulation and Galerkin approximation lead to the least-squares temporal difference
 691 (LSTD) estimator; this connection between Galerkin methods and LSTD has been noted in past work
 692 by Yu and Bertsekas [YB10]. In this paper, our goal is to understand the weak formulation (25) in a
 693 broader sense for general test and predictor classes.

694 **A.2 Comparison with Weight Learning Methods**

695 The work closest to ours is [JH20]. They also use an auxiliary weight function class, which is
 696 comparable to our test class. However, the test class is used in different ways; we compare them
 697 in this section at the population level.⁵ Let us assume that weak realizability holds and that \mathcal{F} is
 698 symmetric, i.e., if $f \in \mathcal{F}$ then $-f \in \mathcal{F}$ as well. At the population level, our program seeks to solve

$$\sup_{Q \in \mathcal{Q}^\pi} \mathbb{E}_{s \sim \nu_{\text{start}}} Q(s, \pi) \quad \text{s.t.} \quad \sup_{f \in \mathcal{F}} \langle f, \mathcal{B}^\pi Q \rangle_\mu = 0, \quad (26)$$

699 which is equivalent for any $w \in \mathcal{F}$ to

$$\sup_{Q \in \mathcal{Q}^\pi} \mathbb{E}_{s \sim \nu_{\text{start}}} Q(s, \pi) - \frac{1}{1 - \gamma} \langle w, \mathcal{B}^\pi Q \rangle_\mu \quad \text{s.t.} \quad \sup_{f \in \mathcal{F}} \langle f, \mathcal{B}^\pi Q \rangle_\mu = 0.$$

⁵The empirical estimator in [JH20] does not take into account the ‘alignment’ of each weight function with respect to the dataset, which we do through self-normalization and regularization in the construction of the empirical estimator. This precludes obtaining the same type of strong finite time guarantees that we are able to derive here.

700 Removing the constraints leads to the upper bound

$$\sup_{Q \in \mathcal{Q}^\pi} \mathbb{E}_{s \sim \nu_{\text{start}}} Q(s, \pi) - \frac{1}{1 - \gamma} \langle w, \mathcal{B}^\pi Q \rangle_\mu.$$

701 Since this is a valid upper bound for any $w \in \mathcal{F}$, minimizing over w must still yield an upper bound,
702 which reads

$$\inf_{w \in \mathcal{F}} \sup_{Q \in \mathcal{Q}^\pi} \mathbb{E}_{s \sim \nu_{\text{start}}} Q(s, \pi) - \frac{1}{1 - \gamma} \langle w, \mathcal{B}^\pi Q \rangle_\mu.$$

703 This is the population program for “weight learning”, as described in [JH20]. It follows that Bellman
704 residual orthogonalization always produces tighter confidence intervals than “weight learning” at the
705 population level.

706 Another interesting comparison is with “value learning”, also described in [JH20]. In this case,
707 assuming symmetric \mathcal{F} , we can equivalently express the population program (26) using a Lagrange
708 multiplier as follows

$$\sup_{Q \in \mathcal{Q}^\pi} \mathbb{E}_{s \sim \nu_{\text{start}}} Q(s, \pi) - \sup_{\lambda \geq 0, f \in \mathcal{F}} \lambda \langle f, \mathcal{B}^\pi Q \rangle_\mu. \quad (27)$$

709 Rearranging we obtain

$$\sup_{Q \in \mathcal{Q}^\pi} \inf_{\lambda \geq 0, f \in \mathcal{F}} \mathbb{E}_{s \sim \nu_{\text{start}}} Q(s, \pi) - \lambda \langle f, \mathcal{B}^\pi Q \rangle_\mu.$$

710 The “value learning” program proposed in [JH20] has a similar formulation to ours but differs in
711 two key aspects. The first—and most important—is that [JH20] ignores the Lagrange multiplier;
712 this means “value learning” is not longer associated to a constrained program. While the Lagrange
713 multiplier could be “incorporated” into the test class \mathcal{F} , doing so would cause the entropy of \mathcal{F} to
714 be unbounded. Another point of difference is that “value learning” uses such expression with $\lambda = 1$
715 to derive the confidence interval *lower bound*, while we use it to construct the confidence interval
716 *upper bound*. While this may seem like a contradiction, we notice that the expression is derived using
717 different assumptions: we assume weak realizability of Q , while [JH20] assumes realizability of the
718 density ratios between μ and the discounted occupancy measure π .

719 A.3 Additional Literature

720 Here we summarize some additional literature. The efficiency of off-policy tabular RL has been
721 investigated in the papers [YBW20, YW20, YW21]. For empirical studies on offline RL, see the
722 papers [LTDC19, JGS⁺19, WTN19, ASN20, WNŽ⁺20, SSB⁺20, NDGL20, YQCC21, KHSL21,
723 BGB20, KFTL19, KRNJ20, YTY⁺20].

724 Some of the classical RL algorithm are presented in the papers [Mun03, Mun05, AMS07, ASM08,
725 FSM10, FGSM16]. For a more modern analysis, see [CJ19]. These works generally make
726 additionally assumptions on top of realizability. Alternatively, one can use importance sam-
727 pling [Pre00, TB16, JL16, FCG18]. A more recent idea is to look at the distributions them-
728 selves [LLTZ18, NDK⁺19, XMW19, ZDLS20, ZLW20, YND⁺20, KU19].

729 Offline policy optimization with pessimism has been studied in the papers [LSAB20, RZM⁺21,
730 JYW21, XCJ⁺21, ZWB21, YWDW, US21]. There exists a fairly extensive literature on lower bounds
731 with linear representations, including the two papers [Zan20, WFK20] that concurrently derived the
732 first exponential lower bounds for the offline setting, and [FKSLX21] proves that realizability and
733 coverage alone are insufficient.

734 In the context of off-policy optimization several works have investigated methods that assume only
735 realizability of the optimal policy [XJ20a, XJ20b]. Related work includes the papers [DW20, DJL21,
736 JH20, UHJ20, TFL⁺19, ND20, VJY21, HJD⁺21, ZSU⁺22, UIJ⁺21, CQ22, LTND21]. Among
737 concurrent works, we note [ZHH⁺22].

738 A.4 Definition of Weak Bellman Closure

739 **Definition 1** (Weak Bellman Closure). *The Bellman operator \mathcal{T}^π is weakly closed with respect to*
740 *the triple $(\mathcal{Q}^\pi, \mathcal{F}^\pi, \mu)$ if for any $Q \in \mathcal{Q}^\pi$, there exists a predictor $\mathcal{P}^\pi(Q) \in \mathcal{Q}^\pi$ such that*

$$\langle f, \mathcal{P}^\pi(Q) \rangle_\mu = \langle f, \mathcal{T}^\pi(Q) \rangle_\mu. \quad (28)$$

741 **A.5 Additional results on the concentrability coefficients**

742 **A.5.1 Testing with the identity function**

743 Suppose that the identity function $\mathbb{1}$ belongs to the test class. Doing so amounts to requiring that the
 744 Bellman error is controlled in an average sense over all the data. When this choice is made, we can
 745 derive some generic upper bounds on K^π , which we state and prove here:

746 **Lemma 4.** *If $\mathbb{1} \in \mathcal{F}^\pi$, then we have the upper bounds*

$$K^\pi \stackrel{(i)}{\leq} \frac{\max_{Q \in \mathcal{C}_n^\pi} |\mathbb{E}_\pi \mathcal{B}^\pi Q|^2}{\max_{Q \in \mathcal{C}_n^\pi} |\mathbb{E}_\mu \mathcal{B}^\pi Q|^2} \stackrel{(ii)}{\leq} K_*^\pi \stackrel{def}{=} \max_{Q \in \mathcal{C}_n^\pi} \frac{|\mathbb{E}_\pi \mathcal{B}^\pi Q|^2}{|\mathbb{E}_\mu \mathcal{B}^\pi Q|^2}. \quad (29)$$

747 *Proof.* Since $\mathbb{1} \in \mathcal{F}$, the definition of \mathcal{C}_n^π implies that

$$\max_{Q \in \mathcal{C}_n^\pi} |\mathbb{E}_\mu \mathcal{B}^\pi Q|^2 \leq (\|\mathbb{1}\|_\mu^2 + \lambda) \frac{\rho}{n} = (1 + \lambda) \frac{\rho}{n}.$$

748 The upper bound (i) then follows from the definition of K^π . The upper bound (ii) follows since the
 749 right hand side is the maximum ratio. \square

750 Note that large values of K_*^π can arise when there exist Q -functions in the set \mathcal{C}_n^π that have low
 751 average Bellman error under the data-generating distribution μ , but relatively large values under π . Of
 752 course, the likelihood of such unfavorable choices of Q is reduced when we use a larger test function
 753 class, which then reduces the size of \mathcal{C}_n^π . However, we pay a price in choosing a larger test function
 754 class, since the choice (40b) of the radius ρ needed for Theorem 3 depends on its complexity.

755 **A.5.2 Mixture distributions**

756 Now suppose that the dataset consists of a collection of trajectories collected by different protocols.
 757 More precisely, for each $j = 1, \dots, m$, let μ_j be a particular protocol for generating a trajectory.
 758 Suppose that we generate data by first sampling a random index $J \in [m]$ according to a probability
 759 distribution $\{p_j\}_{j=1}^m$, and conditioned $J = j$, we sample (s, a, o) according to μ_j . The resulting data
 760 follows a mixture distribution, where we set $o = j$ to tag the protocol used to generate the data. To
 761 be clear, for each sample $i = 1, \dots, n$, we sample J as described, and then draw a single sample
 762 $(s, a, o) \sim \mu_j$.

763 Following the intuition given in the previous section, it is natural to include test functions that code
 764 for the protocol—that is, the binary-indicator functions

$$f_j(s, a, o) = \begin{cases} 1 & \text{if } o = j \\ 0 & \text{otherwise.} \end{cases} \quad (30)$$

765 This test function, when included in the weak formulation, enforces the Bellman evaluation equations
 766 for the policy $\pi \in \Pi$ under consideration along the distribution induced by each data-generating
 767 policy μ_j .

768 **Lemma 5** (Mixture Policy Concentrability). *Suppose that μ is an m -component mixture, and that
 769 the indicator functions $\{f_j\}_{j=1}^m$ are included in the test class. Then we have the upper bounds*

$$K^\pi \stackrel{(i)}{\leq} \frac{1 + m\lambda}{1 + \lambda} \frac{\max_{Q \in \mathcal{C}_n^\pi} |\mathbb{E}_\pi \mathcal{B}^\pi Q|^2}{\max_{Q \in \mathcal{C}_n^\pi} \sum_{j=1}^m p_j^2 |\mathbb{E}_{\mu_j} \mathcal{B}^\pi Q|^2} \stackrel{(ii)}{\leq} \frac{1 + m\lambda}{1 + \lambda} \max_{Q \in \mathcal{C}_n^\pi} \left\{ \frac{|\mathbb{E}_\pi \mathcal{B}^\pi Q|^2}{\sum_{j=1}^m p_j^2 |\mathbb{E}_{\mu_j} \mathcal{B}^\pi Q|^2} \right\}. \quad (31)$$

770 *Proof.* From the definition of K^π , it suffices to show that

$$\max_{Q \in \mathcal{C}_n^\pi} \sum_{j=1}^m p_j^2 |\mathbb{E}_{\mu_j} \mathcal{B}^\pi Q|^2 \leq \frac{\rho}{n} (1 + m\lambda).$$

771 A direct calculation yields $\langle f_j, \mathcal{B}^\pi Q \rangle_\mu = \mathbb{E}_\mu \mathbb{I}\{o = j\} \mathcal{B}^\pi Q = p_j \mathbb{E}_{\mu_j} \mathcal{B}^\pi Q$. Moreover, since each f_j
 772 belongs to the test class by assumption, we have the upper bound $|p_j \mathbb{E}_{\mu_j} \mathcal{B}^\pi Q| \leq \sqrt{\frac{\rho}{n}} \sqrt{\|f_j\|_\mu^2 + \lambda}$.

773 Squaring each term and summing over the constraints yields

$$\sum_{j=1}^m p_j^2 [\mathbb{E}_{\mu_j} \mathcal{B}^\pi Q]^2 \leq \frac{\rho}{n} \sum_{j=1}^m (\|f_j\|_\mu^2 + \lambda) = \frac{\rho}{n} (1 + m\lambda),$$

774 where the final equality follows since $\sum_{j=1}^m \|f_j\|_\mu^2 = 1$. \square

775 As shown by the upper bound, the off-policy coefficient K^π provides a measure of how the squared-
 776 averaged Bellman errors along the policies $\{\mu_j\}_{j=1}^m$, weighted by their probabilities $\{p_j\}_{j=1}^m$, trans-
 777 fers to the evaluation policy π . Note that the regularization parameter λ decays as a function of the
 778 sample size—e.g., as $1/n$ in Theorem 3—the factor $(1 + m\lambda)/(1 + \lambda)$ approaches one as n increases
 779 (for a fixed number m of mixture components).

780 A.5.3 Bellman Rank for off-policy evaluation

781 In this section, we show how more refined bounds can be obtained when—in addition to a mixture
 782 condition—additional structure is imposed on the problem. In particular, we consider a notion similar
 783 to that of Bellman rank [JKA⁺17], but suitably adapted⁶ to the off-policy setting.

784 Given a policy class $\tilde{\Pi}$ and a predictor class $\tilde{\mathcal{Q}}$, we say that it has Bellman rank is d if there exist two
 785 maps $\nu : \tilde{\Pi} \rightarrow \mathbb{R}^d$ and $\xi : \tilde{\mathcal{Q}} \rightarrow \mathbb{R}^d$ such that

$$\mathbb{E}_\pi \mathcal{B}^\pi Q = \langle \nu_\pi, \xi_Q \rangle_{\mathbb{R}^d}, \quad \text{for all } \pi \in \tilde{\Pi} \text{ and } Q \in \tilde{\mathcal{Q}}. \quad (32)$$

786 In words, the average Bellman error of any predictor Q along any given policy π can be expressed
 787 as the Euclidean inner product between two d -dimensional vectors, one for the policy and one for
 788 the predictor. As in the previous section, we assume that the data is generated by a mixture of m
 789 different distributions (or equivalently policies) $\{\mu_j\}_{j=1}^m$. In the off-policy setting, we require that
 790 the policy class $\tilde{\Pi}$ contains all of these policies as well as the target policy—viz. $\{\mu_j\} \cup \{\pi\} \subseteq \tilde{\Pi}$.
 791 Moreover, the predictor class $\tilde{\mathcal{Q}}$ should contain the predictor class for the target policy, i.e., $\mathcal{Q}^\pi \subseteq \tilde{\mathcal{Q}}$.
 792 We also assume weak realizability for this discussion.

793 Our result depends on a positive semidefinite matrix determined by the mixture weights $\{p_j\}_{j=1}^m$
 794 along with the embeddings $\{\nu_{\mu_j}\}_{j=1}^m$ of the associated policies that generated the data. In particular,
 795 we define

$$\Sigma_\nu = \sum_{j=1}^m p_j^2 \nu_{\mu_j} \nu_{\mu_j}^\top.$$

796 Assuming that this matrix is positive definite,⁷ we define the norm $\|u\|_{\Sigma_\nu^{-1}} = \sqrt{u^\top (\Sigma_\nu)^{-1} u}$. With
 797 this notation, we have the following bound.

798 **Lemma 6** (Concentrability with Bellman Rank). *For a mixture data-generation process and under*
 799 *the Bellman rank condition (32), we have the upper bound*

$$K^\pi \leq \frac{1 + m\lambda}{1 + \lambda} \|\nu_\pi\|_{\Sigma_\nu^{-1}}^2, \quad (33)$$

800 *Proof.* Our proof exploits the upper bound (ii) from the claim (31) in Lemma 5. We first evaluate and
 801 redefine the ratio in this upper bound. Weak realizability coupled with the Bellman rank condition (32)
 802 implies that there exists some Q_\star^π such that

$$\begin{aligned} 0 &= \langle f_j, \mathcal{B}^\pi Q_\star^\pi \rangle_\mu = p_j \mathbb{E}_{\mu_j} \mathcal{B}^\pi Q_\star^\pi = p_j \langle \nu_{\mu_j}, \xi_{Q_\star^\pi} \rangle, & \text{for all } j = 1, \dots, m, \text{ and} \\ 0 &= \langle \mathbf{1}, \mathcal{B}^\pi Q_\star^\pi \rangle_\pi = \mathbb{E}_\pi \mathcal{B}^\pi Q_\star^\pi = \langle \nu_\pi, \xi_{Q_\star^\pi} \rangle. \end{aligned}$$

⁶The original definition essentially takes $\tilde{\Pi}$ as the set of all greedy policies with respect to $\tilde{\mathcal{Q}}$. Since a dataset need not originate from greedy policies, the definition of Bellman rank is adapted in a natural way.

⁷If not, one can prove a result for a suitably regularized version.

803 Therefore, we have the equivalences $\mathbb{E}_{\mu_j} \mathcal{B}^\pi Q = \langle \nu_{\mu_j}, (\xi_Q - \xi_{Q_\star^\pi}) \rangle$ for all $j = 1, \dots, m$, as well as
804 $\mathbb{E}_\pi \mathcal{B}^\pi Q = \langle \nu_\pi, (\xi_Q - \xi_{Q_\star^\pi}) \rangle$. Introducing the shorthand $\Delta_Q = \xi_Q - \xi_{Q_\star^\pi}$, we can bound the ratio
805 as follows

$$\begin{aligned} \sup_{Q \in \mathcal{C}_\pi^\pi} \left\{ \frac{(\langle \nu_\pi, \Delta_Q \rangle)^2}{\sum_{j=1}^m p_j^2 (\langle \nu_{\mu_j}, \Delta_Q \rangle)^2} \right\} &= \sup_{Q \in \mathcal{C}_\pi^\pi} \left\{ \frac{(\langle \nu_\pi, \Delta_Q \rangle)^2}{\Delta_Q^\top \left(\sum_{j=1}^m p_j^2 \nu_{\mu_j} \nu_{\mu_j}^\top \right) \Delta_Q} \right\} \\ &= \sup_{Q \in \mathcal{C}_\pi^\pi} \left\{ \frac{(\langle \nu_\pi, \Sigma_\nu^{-\frac{1}{2}} \tilde{\Delta}_Q \rangle)^2}{\|\tilde{\Delta}_Q\|_2^2} \right\} \quad \text{where } \tilde{\Delta}_Q = \Sigma_\nu^{\frac{1}{2}} \Delta_Q \\ &\leq \|\nu_\pi\|_{\Sigma_\nu^{-1}}^2, \end{aligned}$$

806 where the final step follows from the Cauchy–Schwarz inequality. \square

807 Thus, when performing off-policy evaluation with a mixture distribution under the Bellman rank
808 condition, the coefficient K^π is bounded by the alignment between the target policy π and the
809 data-generating distribution μ , as measured in the the embedded space guaranteed by the Bellman
810 rank condition. The structure of this upper bound is similar to a result that we derive in the sequel for
811 linear approximation under Bellman closure (see Proposition 3).

812 A.6 Further comments on the prediction error test space

813 A few comments on the bound in Lemma 1: as in our previous results, the pre-factor $\frac{\|\epsilon\|_\mu^2 + \lambda}{\|\mathbb{1}\|_\pi^2 + \lambda}$
814 serves as a normalization factor. Disregarding this leading term, the second ratio measures how the
815 prediction error $\epsilon = Q - Q_\star^\pi$ along μ transfers to π , as measured via the operator $\mathcal{L} - \gamma \mathbb{P}^\pi$. This
816 interaction is complex, since it includes the *bootstrapping term* $-\gamma \mathbb{P}^\pi$. (Notably, such a term is not
817 present for standard prediction or bandit problems, in which case $\gamma = 0$.) This term reflects the
818 dynamics intrinsic to reinforcement learning, and plays a key role in proving “hard” lower bounds for
819 offline RL (e.g., see the work [Zan20]).

820 Observe that the bound in Lemma 1 requires only weak realizability, and thus it always applies. This
821 fact is significant in light of a recent lower bound [FKSLX21], showing that without Bellman closure,
822 off-policy learning is challenging even under strong concentrability assumption (such as bounds on
823 density ratios). Lemma 1 gives a sufficient condition without Bellman closure, but with a different
824 measure that accounts for bootstrapping.

825

826 If, in fact, (weak) Bellman closure holds, then Lemma 1 takes the following simplified form:

827 **Lemma 7** (OPC coefficient under Bellman closure). *If $\mathcal{E}^\pi \subseteq \mathcal{F}^\pi$ and weak Bellman closure holds,*
828 *then*

$$K^\pi \leq \max_{\epsilon \in \mathcal{E}^\pi} \left\{ \frac{\|\epsilon\|_\mu^2 + \lambda}{1 + \lambda} \cdot \frac{\langle \mathbb{1}, \epsilon \rangle_\pi^2}{\langle \epsilon, \epsilon \rangle_\mu} \right\} \leq \max_{\epsilon \in \mathcal{E}^\pi} \left\{ \frac{\|\epsilon\|_\pi^2}{\|\epsilon\|_\mu^2} \right\}.$$

829 See Appendix D.3 for the proof.

830

831 In such case, the concentrability measures the increase in the discrepancy $Q - Q'$ of the feasible
832 predictors when moving from the dataset distribution μ to the distribution of the target policy π . In
833 Section 4.3, we give another bound under weak Bellman closure, and thereby recover a recent result
834 due to Xie et al. [XCJ⁺21]. Finally, in Section 5, we provide some applications of this concentrability
835 factor to the linear setting.

836 A.7 From Importance Sampling to Bellman Closure

837 Let us show an application of Lemma 3 on an example with just two test spaces. Suppose that we
838 suspect that Bellman closure holds, but rather than committing to such assumption, we wish to fall
839 back to an importance sampling estimator if Bellman closure does not hold.

840 In order to streamline the presentation of the idea, let us introduce the following setup. Let π^b be a
841 behavioral policy that generates the dataset, i.e., such that each state-action (s, a) in the dataset is

842 sampled from its discounted state distribution d_{π^b} . Next, let the identifier o contain the trajectory
843 from ν_{start} up to the state-action pair (s, a) recorded in the dataset. That is, each tuple (s, a, r, s^+, o)
844 in the dataset \mathcal{D} is such that $(s, a) \sim d_{\pi^b}$ and o contains the trajectory up to (s, a) .

845 We now define the test spaces. The first one is denoted with $\mathcal{F}_\pi^{\text{IS}}$ and leverages importance sampling.
846 It contains a single test function defined as the importance sampling estimator

$$\mathcal{F}_\pi^{\text{IS}} = \{f_\pi\}, \quad \text{where } f_\pi(s, a, o) = \frac{1}{b_\pi} \prod_{(s_h, a_h) \in o} \frac{\pi(a_h | s_h)}{\pi^b(a_h | s_h)}. \quad (34)$$

847 The above product is over the random trajectory contained in the identifier o . The normalization
848 factor $b_\pi \in \mathbb{R}$ is connected to the maximum range of the importance sampling estimator, and ensures
849 that $\sup_{(s, a, o)} f_\pi(s, a, o) \leq 1$. The second test space is the prediction error test space \mathcal{E}^π defined in
850 Section 4.2.

851 With this choice, let us define three concentrability coefficients. $K_{(1)}^\pi$ arises from importance sampling,
852 $K_{(2)}^\pi$ from the prediction error test space when Bellman closure holds and $K_{(3)}^\pi$ from the prediction
853 error test space when just weak realizability holds. They are defined as

$$K_{(1)}^\pi \leq \sqrt{b_\pi \frac{(1 + \lambda b_\pi)}{1 + \lambda}} \quad K_{(2)}^\pi \leq \max_{\epsilon \in \mathcal{E}_\pi^*} \frac{\langle \mathbf{1}, (\mathcal{I} - \gamma \mathbb{P}^\pi) \epsilon \rangle_\pi^2}{\langle \epsilon, (\mathcal{I} - \gamma \mathbb{P}^\pi) \epsilon \rangle_\mu^2} \times \frac{\|\epsilon\|_\mu^2 + \lambda}{\|\mathbf{1}\|_\pi^2 + \lambda}, \quad K_{(3)}^\pi \leq c_1 \frac{\|\mathcal{B}^\pi Q\|_\pi^2}{\|\mathcal{B}^\pi Q\|_\mu^2}.$$

854 **Lemma 8** (From Importance Sampling to Bellman Closure). *The choice $\mathcal{F}^\pi = \mathcal{F}_\pi^{\text{IS}} \cup \mathcal{E}^\pi$ for all $\pi \in$
855 Π ensures that with probability at least $1 - \delta$, the oracle inequality (9) holds with $K^\pi \leq$
856 $\min\{K_{(1)}^\pi, K_{(2)}^\pi, K_{(3)}^\pi\}$ if weak Bellman closure holds and $K^\pi \leq \min\{K_{(1)}^\pi, K_{(2)}^\pi\}$ otherwise.*

857 *Proof.* Let us calculate the off-policy cost coefficient associated with $\mathcal{F}_\pi^{\text{IS}}$. The unbiasedness of the
858 importance sampling estimator gives us the following population constraint (here $\mu = d_{\pi^b}$)

$$|\langle f_\pi, \mathcal{B}^\pi Q \rangle_\mu| = |\mathbb{E}_\mu f_\pi \mathcal{B}^\pi Q| = \frac{1}{b_\pi} |\mathbb{E}_\pi \mathcal{B}^\pi Q| = \frac{1}{b_\pi} |\langle \mathbf{1}, \mathcal{B}^\pi Q \rangle_\pi| \leq \frac{L}{\sqrt{n}} \sqrt{\|f_\pi\|_2^2 + \lambda}$$

859 The norm of the test function reads (notice that μ generates (s, a, o) here)

$$\|f_\pi\|_\mu^2 = \mathbb{E}_\mu f_\pi^2 = \frac{1}{b_\pi^2} \mathbb{E}_\mu \left[\prod_{(s_h, a_h) \in o} \frac{\pi(a_h | s_h)}{\pi^b(a_h | s_h)} \right]^2 = \frac{1}{b_\pi^2} \mathbb{E}_\pi \left[\prod_{(s_h, a_h) \in o} \frac{\pi(a_h | s_h)}{\pi^b(a_h | s_h)} \right] \leq \frac{1}{b_\pi}.$$

860 Together with the prior display, we obtain

$$\frac{\langle \mathbf{1}, \mathcal{B}^\pi Q \rangle_\pi^2}{b_\pi^2 (\|f_\pi\|_2^2 + \lambda)} \leq \frac{\rho}{n}.$$

861 The resulting concentrability coefficient is therefore

$$K^\pi \leq \max_{Q \in \mathcal{C}_\pi^n} \frac{\langle \mathbf{1}, \mathcal{B}^\pi Q \rangle_\pi^2}{1 + \lambda} \times \frac{n}{\rho} \leq \max_{Q \in \mathcal{C}_\pi^n} \frac{\langle \mathbf{1}, \mathcal{B}^\pi Q \rangle_\pi^2}{1 + \lambda} \times \frac{b_\pi^2 (\|f_\pi\|_2^2 + \lambda)}{\langle \mathbf{1}, \mathcal{B}^\pi Q \rangle_\pi^2} \leq b_\pi \frac{(1 + \lambda b_\pi)}{1 + \lambda}.$$

862 Chaining the above result with Lemmas 1 and 2, using Lemma 3 and plugging back into Theorem 3
863 yields the thesis. \square

864 A.8 Implementation for Off-Policy Predictions

865 In this section, we describe a computationally efficient way in which to compute the upper/lower
866 estimates (5). Given a finite set of $n_{\mathcal{F}}$ test functions, it involves solving a quadratic program with
867 $2n_{\mathcal{F}} + 1$ constraints.

868 Let us first work out a concise description of the constraints defining membership in $\widehat{\mathcal{C}}_n^\pi$. Introduce the
869 shorthand $n_f \stackrel{\text{def}}{=} \|f_j\|_n^2 + \lambda$. We then define the empirical average feature vector $\widehat{\phi}_f$, the empirical

870 average reward \hat{r}_f , and the average next-state feature vector $\hat{\phi}_f^{+\pi}$ as

$$\begin{aligned}\hat{\phi}_f &= \frac{1}{\sqrt{n_f}} \sum_{(s,a,r,s^+) \in \mathcal{D}} f(s,a)\phi(s,a), & \hat{r}_f &= \frac{1}{\sqrt{n_f}} \sum_{(s,a,r,s^+) \in \mathcal{D}} f(s,a)r, \\ \hat{\phi}_f^{+\pi} &= \frac{1}{\sqrt{n_f}} \sum_{(s,a,r,s^+) \in \mathcal{D}} f(s,a)\phi(s^+, \pi).\end{aligned}$$

871 In terms of this notation, each empirical constraint defining $\hat{\mathcal{C}}_n^\pi$ can be written in the more compact
872 form

$$\frac{|\langle f, \delta^\pi Q \rangle_n|}{\sqrt{n_f}} = \left| \langle \hat{\phi}_f - \gamma \hat{\phi}_f^{+\pi}, w \rangle - \hat{r}_f \right| \leq \sqrt{\frac{\rho}{n}}.$$

873 Then the set of empirical constraints can be written as a set of constraints linear in the critic parameter
874 w coupled with the assumed regularity bound on w

$$\hat{\mathcal{C}}_n^\pi = \left\{ w \in \mathbb{R}^d \mid \|w\|_2 \leq 1, \quad \text{and} \quad -\sqrt{\frac{\rho}{n}} \leq \langle \hat{\phi}_f - \gamma \hat{\phi}_f^{+\pi}, w \rangle - \hat{r}_f \leq \sqrt{\frac{\rho}{n}} \quad \text{for all } f \in \mathcal{F}^\pi \right\}. \quad (35)$$

875 Thus, the estimates \hat{V}_{\min}^π (respectively \hat{V}_{\max}^π) can be computed by minimizing (respectively maximiz-
876 ing) the linear objective function $w \mapsto \langle [\mathbb{E}_{s \sim \nu_{\text{start}}} \mathbb{E}_{a \sim \pi} \phi(s,a)], w \rangle$ subject to the $2n_{\mathcal{F}} + 1$ constraints
877 in equation (35). Therefore, the estimates can be computed in polynomial time for any test function
878 with a cardinality that grows polynomially in the problem parameters.

879 A.9 Discussion of Linear Approximate Optimization

880 Here we discuss the presence of the supremum over policies in the coefficient $K_{(1)}^{\tilde{\pi}}$ from equation (23).
881 In particular, it arises because our actor-critic method iteratively approximates the maximum in the
882 max-min estimate (6) using a gradient-based scheme. The ability of a gradient-based method to make
883 progress is related to the estimation accuracy of the gradient, which is the Q estimates of the actor's
884 current policy π_t ; more specifically, the gradient is the Q function parameter w_t . In the general case,
885 the estimation error of the gradient w_t depends on the policy under consideration through the matrix
886 $\Sigma_{\lambda, \text{Boot}}^{+\pi_t}$, while it is independent in the special case of Bellman closure (as it depends on just Σ). As the
887 actor's policies are random, this yields the introduction of a $\sup_{\pi \in \Pi}$ in the general bound. Notice the
888 method still competes with the best comparator $\tilde{\pi}$ by measuring the errors along the distribution of the
889 comparator (through the operator $\mathbb{E}_{\tilde{\pi}}$). To be clear, $\sup_{\pi \in \Pi}$ may not arise with approximate solution
890 methods that do not rely only on the gradient to make progress (such as second-order methods); we
891 leave this for future research. Reassuringly, when Bellman closure, the approximate solution method
892 recovers the standard guarantees established in the paper [ZWB21].

893 **B General Guarantees**

894 **B.1 A deterministic guarantee**

895 We begin our analysis stating a deterministic set of sufficient conditions for our estimators to satisfy
 896 the guarantees (8) and (9). This formulation is useful, because it reveals the structural conditions
 897 that underlie success of our estimators, and in particular the connection to weak realizability. In
 898 Section B.2, we exploit this deterministic result to show that, under a fairly general sampling model,
 899 our estimators enjoy these guarantees with high probability.

900 In the previous section, we introduced the population level set \mathcal{C}_n^π that arises in the statement of our
 901 guarantees. Also central in our analysis is the infinite data limit of this set. More specifically, for
 902 any fixed (ρ, λ) , if we take the limit $n \rightarrow \infty$, then \mathcal{C}_n^π reduces to the set of all solutions to the weak
 903 formulation (25)—that is

$$\mathcal{C}_\infty^\pi(\mathcal{F}^\pi) = \{Q \in \mathcal{Q}^\pi \mid \langle f, \mathcal{B}^\pi Q \rangle_\mu = 0 \text{ for all } f \in \mathcal{F}^\pi\}. \quad (36)$$

904 As before, we omit the dependence on the test function class \mathcal{F}^π when it is clear from context. By
 905 construction, we have the inclusion $\mathcal{C}_\infty^\pi(\mathcal{F}^\pi) \subseteq \mathcal{C}_n^\pi(4\rho, \lambda; \mathcal{F}^\pi)$ for any non-negative pair (ρ, λ) .

906 Our first set of guarantees hold when the random set $\widehat{\mathcal{C}}_n^\pi$ satisfies the *sandwich relation*

$$\mathcal{C}_\infty^\pi(\mathcal{F}^\pi) \subseteq \widehat{\mathcal{C}}_n^\pi(\rho, \lambda; \mathcal{F}^\pi) \subseteq \mathcal{C}_n^\pi(4\rho, \lambda; \mathcal{F}^\pi) \quad (37)$$

907 To provide intuition as to why this sandwich condition is natural, observe that it has two important
 908 implications:

- 909 (a) Recalling the definition of weak realizability (1), the weak solution Q_*^π belongs to the
 910 empirical constraint set $\widehat{\mathcal{C}}_n^\pi$ for any choice of test function space. This important property
 911 follows because Q_*^π must satisfy the constraints (25), and thus it belongs to $\mathcal{C}_\infty^\pi \subseteq \widehat{\mathcal{C}}_n^\pi$.
- 912 (b) All solutions in $\widehat{\mathcal{C}}_n^\pi$ also belong to \mathcal{C}_n^π , which means they approximately satisfy the weak
 913 Bellman equations in a way quantified by \mathcal{C}_n^π .

914 By leveraging these facts in the appropriate way, we can establish the following guarantee:

915

916 **Proposition 4.** *The following two statements hold.*

- 917 (a) *Policy evaluation:* *If the set $\widehat{\mathcal{C}}_n^\pi$ satisfies the sandwich relation (37), then the estimates*
 918 *$(\widehat{V}_{\min}^\pi, \widehat{V}_{\max}^\pi)$ satisfy the width bound (8b). If, in addition, weak Bellman realizability for π is*
 919 *assumed, then the coverage (8a) condition holds.*
- 920 (b) *Policy optimization:* *If the sandwich relation (37) and weak Bellman realizability hold for*
 921 *all $\pi \in \Pi$, then any max-min (6) optimal policy $\tilde{\pi}$ satisfies the oracle inequality (9).*

922 See Section C.1 for the proof of this claim.

923

924 In summary, Proposition 4 ensures that when weak realizability is in force, then the sandwich
 925 relation (37) is a sufficient condition for both the policy evaluation (8) and optimization (9) guarantees
 926 to hold. Accordingly, the next phase of our analysis focuses on deriving sufficient conditions for the
 927 sandwich relation to hold with high probability.

928 **B.2 Some high-probability guarantees**

929 As stated, Proposition 4 is a “meta-result”, in that it applies to any choice of set $\widehat{\mathcal{C}}_n^\pi \equiv \widehat{\mathcal{C}}_n^\pi(\rho, \lambda; \mathcal{F}^\pi)$
 930 for which the sandwich relation (37) holds. In order to obtain a more concrete guarantee, we need to
 931 impose assumptions on the way in which the dataset was generated, and concrete choices of (ρ, λ)
 932 that suffice to ensure that the associated sandwich relation (37) holds with high probability. These
 933 tasks are the focus of this section.

934 **B.2.1 A model for data generation**

935 Let us begin by describing a fairly general model for data-generation. Any sample takes the form
 936 $z \stackrel{def}{=} (s, a, r, s^+, o)$, where the five components are defined as follows:

- 937 • the pair (s, a) index the current state and action.
- 938 • the random variable r is a noisy observation of the mean reward.
- 939 • the random state s^+ is the next-state sample, drawn according to the transition $\mathbb{P}(s, a)$.
- 940 • the variable o is an optional identifier.

941 As one example of the use of an identifier variable, if samples might be generated by one of two
 942 possible policies—say π_1 and π_2 —the identifier can take values in the set $\{1, 2\}$ to indicate which
 943 policy was used for a particular sample.

944

945 Overall, we observe a dataset $\mathcal{D} = \{z_i\}_{i=1}^n$ of n such quintuples. In the simplest of possible settings,
 946 each triple (s, a, o) is drawn i.i.d. from some fixed distribution μ , and the noisy reward r_i is an
 947 unbiased estimate of the mean reward function $R(s_i, a_i)$. In this case, our dataset consists of n i.i.d.
 948 quintuples. More generally, we would like to accommodate richer sampling models in which the
 949 sample $z_i = (s_i, a_i, o_i, r_i, s_i^+)$ at a given time i is allowed to depend on past samples. In order to
 950 specify such dependence in a precise way, define the nested sequence of sigma-fields

$$\mathcal{F}_1 = \emptyset, \quad \text{and} \quad \mathcal{F}_i \stackrel{def}{=} \sigma\left(\{z_j\}_{j=1}^{i-1}\right) \quad \text{for } i = 2, \dots, n. \quad (38)$$

951 In terms of this filtration, we make the following definition:

952 **Assumption 3** (Adapted dataset). *An adapted dataset is a collection $\mathcal{D} = \{z_i\}_{i=1}^n$ such that for each*
 953 *$i = 1, \dots, n$:*

- 954 • *There is a conditional distribution μ_i such that $(s_i, a_i, o_i) \sim \mu_i(\cdot \mid \mathcal{F}_i)$.*
- 955 • *Conditioned on (s_i, a_i, o_i) , we observe a noisy reward $r_i = r(s_i, a_i) + \eta_i$ with $\mathbb{E}[\eta_i \mid \mathcal{F}_i] = 0$,*
 956 *and $|r_i| \leq 1$.*
- 957 • *Conditioned on (s_i, a_i, o_i) , the next state s_i^+ is generated according to $\mathbb{P}(s_i, a_i)$.*

958 Under this assumption, we can define the (possibly) random reference measure

$$\mu(s, a, o) \stackrel{def}{=} \frac{1}{n} \sum_{i=1}^n \mu_i(s, a, o \mid \mathcal{F}_i). \quad (39)$$

959 In words, it corresponds to the distribution induced by first drawing a time index $i \in \{1, \dots, n\}$
 960 uniformly at random, and then sampling a triple (s, a, o) from the conditional distribution $\mu_i(\cdot \mid \mathcal{F}_i)$.

961 **B.2.2 A general guarantee**

962 Recall that there are three function classes that underlie our method: the test function class \mathcal{F} , the
 963 policy class Π , and the Q -function class \mathcal{Q} . In this section, we state a general guarantee (Theorem 3)
 964 that involves the metric entropies of these sets. In Section B.2.3, we provide corollaries of this
 965 guarantee for specific function classes.

966 In more detail, we equip the test function class and the Q -function class with the usual sup-norm

$$\|f - \tilde{f}\|_\infty \stackrel{def}{=} \sup_{(s,a,o)} |f(s, a, o) - \tilde{f}(s, a, o)|, \quad \text{and} \quad \|Q - \tilde{Q}\|_\infty \stackrel{def}{=} \sup_{(s,a)} |Q(s, a) - \tilde{Q}(s, a)|,$$

967 and the policy class with the sup-TV norm

$$\|\pi - \tilde{\pi}\|_{\infty,1} \stackrel{def}{=} \sup_s \|\pi(\cdot \mid s) - \tilde{\pi}(\cdot \mid s)\|_1 = \sup_s \sum_a |\pi(a \mid s) - \tilde{\pi}(a \mid s)|.$$

968 For a given $\epsilon > 0$, we let $\mathcal{N}_\epsilon(\mathcal{F})$, $\mathcal{N}_\epsilon(\mathcal{Q})$, and $\mathcal{N}_\epsilon(\Pi)$ denote the ϵ -covering numbers of each of these
 969 function classes in the given norms. Given these covering numbers, a tolerance parameter $\delta \in (0, 1)$

970 and the shorthand $\phi(t) = \max\{t, \sqrt{t}\}$, define the radius function

$$\rho(\epsilon, \delta) \stackrel{\text{def}}{=} n \left\{ \int_{\epsilon^2}^{\epsilon} \phi\left(\frac{\log N_u(\mathcal{F})}{n}\right) du + \frac{\log N_\epsilon(\mathcal{Q})}{n} + \frac{\log N_\epsilon(\Pi)}{n} + \frac{\log(n/\delta)}{n} \right\}. \quad (40a)$$

971 In our theorem, we implement the estimator using a radius $\rho = \rho(\epsilon, \delta)$, where $\epsilon > 0$ is any parameter
972 that satisfies the bound

$$\epsilon^2 \stackrel{(i)}{\leq} \bar{c} \frac{\rho(\epsilon, \delta)}{n}, \quad \text{and} \quad \lambda \stackrel{(i)}{=} 4 \frac{\rho(\epsilon, \delta)}{n}. \quad (40b)$$

973 Here $\bar{c} > 0$ is a suitably chosen but universal constant (whose value is determined in the proof), and
974 we adopt the shorthand $\rho = \rho(\epsilon, \delta)$ in our statement below.

975 **Theorem 3** (High-probability guarantees). *Consider the estimates implemented using triple $(\Pi, \mathcal{F}, \mathcal{Q})$
976 that is weakly Bellman realizable (Assumption 1); an adapted dataset (Assumption 3); and with the
977 choices (40) for $(\epsilon, \rho, \lambda)$. Then with probability at least $1 - \delta$:*

978 Policy evaluation: *For any $\pi \in \Pi$, the estimates $(\widehat{V}_{\min}^\pi, \widehat{V}_{\max}^\pi)$ specify a confidence interval satisfying
979 the coverage (8a) and width bounds (8b).*

980 Policy optimization: *Any max-min policy (6) $\tilde{\pi}$ satisfies the oracle inequality (9).*

981 See Appendix C.3 for the proof of the claim.

982

983 **Choices of $(\rho, \epsilon, \lambda)$:** Let us provide a few comments about the choices of $(\rho, \epsilon, \lambda)$ from equa-
984 tions (40a) and (40b). The quality of our bounds depends on the size of the constraint set \mathcal{C}_n^π , which
985 is controlled by the constraint level $\sqrt{\frac{\rho}{n}}$. Consequently, our results are tightest when $\rho = \rho(\epsilon, \delta)$ is as
986 small as possible. Note that ρ is an decreasing function of ϵ , so that in order to minimize it, we would
987 like to choose ϵ as large as possible subject to the constraint (40b)(i). Ignoring the entropy integral
988 term in equation (40b) for the moment—see below for some comments on it—these considerations
989 lead to

$$n\epsilon^2 \asymp \log N_\epsilon(\mathcal{F}) + \log N_\epsilon(\mathcal{Q}) + \log N_\epsilon(\Pi). \quad (41)$$

990 This type of relation for the choice of ϵ in non-parametric statistics is well-known (e.g., see Chapters
991 13–15 in the book [Wai19] and references therein). Moreover, setting $\lambda \asymp \epsilon^2$ as in equation (40b)(ii)
992 is often the correct scale of regularization.

993 **Key technical steps in proof:** It is worthwhile making a few comments about the structure of the
994 proof so as to clarify the connections to Proposition 4 along with the weak formulation that underlies
995 our methods. Recall that Proposition 4 requires the empirical $\widehat{\mathcal{C}}_n^\pi$ and population sets \mathcal{C}_n^π to satisfy the
996 sandwich relation (37). In order to prove that this condition holds with high probability, we need to
997 establish uniform control over the family of random variables

$$\frac{|\langle f, \delta^\pi(Q) \rangle_n - \langle f, \mathcal{B}^\pi(Q) \rangle_\mu|}{\sqrt{\|f\|_n^2 + \lambda}}, \quad \text{as indexed by the triple } (f, Q, \pi). \quad (42)$$

998 Note that the differences in the numerator of these variables correspond to moving from the empirical
999 constraints on Q -functions that are enforced using the TD errors, to the population constraints that
1000 involve the Bellman error function.

1001 Uniform control of the family (42), along with the differences $\|f\|_n - \|f\|_\mu$ uniformly over f ,
1002 allows us to relate the empirical and population sets, since the associated constraints are obtained by
1003 shifting between the empirical inner products $\langle \cdot, \cdot \rangle_n$ to the reference inner products $\langle \cdot, \cdot \rangle_\mu$. A simple
1004 discretization argument allows us to control the differences uniformly in (Q, π) , as reflected by the
1005 metric entropies appearing in our definition (40). Deriving uniform bounds over test functions f —due
1006 to the self-normalizing nature of the constraints—requires a more delicate argument. More precisely,
1007 in order to obtain optimal results for non-parametric problems (see Corollary 2 to follow), we need
1008 to localize the empirical process at a scale ϵ , and derive bounds on the localized increments. This
1009 portion of the argument leads to the entropy integral—which is localized to the interval $[\epsilon^2, \epsilon]$ —in
1010 our definition (40a) of the radius function.

1011 **Intuition from the on-policy setting:** In order to gain intuition for the statistical meaning of the
 1012 guarantees in Theorem 3, it is worthwhile understanding the implications in a rather special case—
 1013 namely, the simpler on-policy setting, where the discounted occupation measure induced by the target
 1014 policy π coincides with the dataset distribution μ . Let us consider the case in which the identity
 1015 function $\mathbb{1}$ belongs to the test class \mathcal{F}^π . Under these conditions, for any $Q \in \mathcal{C}_n^\pi$, we can write

$$\max_{Q \in \mathcal{C}_n^\pi} |\mathbb{E}_\pi \mathcal{B}^\pi Q| \stackrel{(i)}{=} \max_{Q \in \mathcal{C}_n^\pi} |\mathbb{E}_\mu \mathcal{B}^\pi Q| \stackrel{(ii)}{\leq} \sqrt{1+\lambda} \sqrt{\frac{\rho}{n}},$$

1016 where equality (i) follows from the on-policy assumption, and step (ii) follows from the definition of
 1017 the set \mathcal{C}_n^π , along with the condition that $\mathbb{1} \in \mathcal{F}^\pi$. Consequently, in the on-policy setting, the width
 1018 bound (8b) ensures that

$$|\widehat{V}_{\min}^\pi - \widehat{V}_{\max}^\pi| \leq 2 \frac{\sqrt{1+\lambda}}{1-\gamma} \sqrt{\frac{\rho}{n}}. \quad (43)$$

1019 In this simple case, we see that the confidence interval scales as $\sqrt{\rho/n}$, where the quantity ρ is related
 1020 to the metric entropy via equation (40b). In the more general off-policy setting, the bound involves
 1021 this term, along with additional terms that reflect the cost of off-policy data. We discuss these issues
 1022 in more detail in Section 4. Before doing so, however, it is useful derive some specific corollaries that
 1023 show the form of ρ under particular assumptions on the underlying function classes, which we now
 1024 do.

1025 B.2.3 Some corollaries

1026 Theorem 3 applies generally to triples of function classes $(\Pi, \mathcal{F}, \mathcal{Q})$, and the statistical error $\sqrt{\frac{\rho(\epsilon, \delta)}{n}}$
 1027 depends on the metric entropies of these function classes via the definition (40a) of $\rho(\epsilon, \delta)$, and the
 1028 choices (40b). As shown in this section, if we make particular assumptions about the metric entropies,
 1029 then we can derive more concrete guarantees.

1030 **Parametric and finite VC classes:** One form of metric entropy, typical for a relatively simple
 1031 function class \mathcal{G} (such as those with finite VC dimension) scales as

$$\log N_\epsilon(\mathcal{G}) \asymp d \log\left(\frac{1}{\epsilon}\right), \quad (44)$$

1032 for some dimensionality parameter d . For instance, bounds of this type hold for linear function
 1033 classes with d parameters, and for finite VC classes (with d proportional to the VC dimension); see
 1034 Chapter 5 of the book [Wai19] for more details.

1035 **Corollary 1.** *Suppose each class of the triple $(\Pi, \mathcal{F}, \mathcal{Q})$ has metric entropy that is at most poly-*
 1036 *nomial (44) of order d . Then for a sample size $n \geq 2d$, the claims of Theorem 3 hold with $\epsilon^2 = d/n$*
 1037 *and*

$$\tilde{\rho}\left(\sqrt{\frac{d}{n}}, \delta\right) \stackrel{def}{=} c \left\{ d \log\left(\frac{n}{d}\right) + \log\left(\frac{n}{\delta}\right) \right\}, \quad (45)$$

1038 where c is a universal constant.

1039 *Proof.* Our strategy is to upper bound the radius ρ from equation (40a), and then show that this
 1040 upper bound $\tilde{\rho}$ satisfies the conditions (40b) for the specified choice of ϵ^2 . We first control the term
 1041 $\log N_\epsilon(\mathcal{F})$. We have

$$\frac{1}{\sqrt{n}} \int_{\epsilon^2}^{\epsilon} \sqrt{\log N_u(\mathcal{F})} du \leq \sqrt{\frac{d}{n}} \int_0^{\epsilon} \sqrt{\log(1/u)} du = \epsilon \sqrt{\frac{d}{n}} \int_0^1 \sqrt{\log(1/(\epsilon t))} dt = c \epsilon \log(1/\epsilon) \sqrt{\frac{d}{n}}.$$

1042 Similarly, we have

$$\frac{1}{n} \int_{\epsilon^2}^{\epsilon} \log N_u(\mathcal{F}) du \leq \frac{d}{n} \left\{ \int_{\epsilon}^1 \log(1/t) dt + \log(1/\epsilon) \right\} \leq c \epsilon \log(1/\epsilon) \frac{d}{n}.$$

1043 Finally, for terms not involving entropy integrals, we have

$$\max \left\{ \frac{\log N_\epsilon(\mathcal{Q})}{n}, \frac{\log N_\epsilon(\Pi)}{n} \right\} \leq c \frac{d}{n} \log(1/\epsilon).$$

1044 Setting $\epsilon^2 = d/n$, we see that the required conditions (40b) hold with the specified choice (45) of
 1045 $\tilde{\rho}$. \square

1046 **Richer function classes:** In the previous section, the metric entropy scaled logarithmically in the
 1047 inverse precision $1/\epsilon$. For other (richer) function classes, the metric entropy exhibits a polynomial
 1048 scaling in the inverse precision, with an exponent $\alpha > 0$ that controls the complexity. More precisely,
 1049 we consider classes of the form

$$\log N_\epsilon(\mathcal{G}) \asymp \left(\frac{1}{\epsilon}\right)^\alpha. \quad (46)$$

1050 For example, the class of Lipschitz functions in dimension d has this type of metric entropy with
 1051 $\alpha = d$. More generally, for Sobolev spaces of functions that have s derivatives (and the s^{th} -derivative
 1052 is Lipschitz), we encounter metric entropies of this type with $\alpha = d/s$. See Chapter 5 of the
 1053 book [Wai19] for further background.

1054 **Corollary 2.** *Suppose that each function class $(\Pi, \mathcal{F}, \mathcal{Q})$ has metric entropy with at most*
 1055 *α -scaling (46) for some $\alpha \in (0, 2)$. Then the claims of Theorem 3 hold with $\epsilon^2 = (1/n)^{\frac{2}{2+\alpha}}$,*
 1056 *and*

$$\tilde{\rho}\left((1/n)^{\frac{1}{2+\alpha}}, \delta\right) = c \left\{ n^{\frac{\alpha}{2+\alpha}} + \log(n/\delta) \right\}. \quad (47)$$

1057 where c is a universal constant.

1058 We note that for standard regression problems over classes with α -metric entropy, the rate $(1/n)^{\frac{2}{2+\alpha}}$
 1059 is well-known to be minimax optimal (e.g., see Chapter 15 in the book [Wai19], as well as references
 1060 therein).

1061 *Proof.* We start by controlling the terms involving entropy integrals. In particular, we have

$$\frac{1}{\sqrt{n}} \int_{\epsilon^2}^{\epsilon} \sqrt{\log N_u(\mathcal{F})} du \leq \frac{c}{\sqrt{n}} u^{1-\frac{\alpha}{2}} \Big|_0^{\epsilon} = \frac{c}{\sqrt{n}} \epsilon^{1-\frac{\alpha}{2}}.$$

1062 Requiring that this term is of order ϵ^2 amounts to enforcing that $\epsilon^{1+\frac{\alpha}{2}} \asymp (1/\sqrt{n})$, or equivalently
 1063 that $\epsilon^2 \asymp (1/n)^{\frac{2}{2+\alpha}}$.

1064 If $\alpha \in (0, 1]$, then the second entropy integral converges and is of lower order. Otherwise, if
 1065 $\alpha \in (1, 2)$, then we have

$$\frac{1}{n} \int_{\epsilon^2}^{\epsilon} \log N_u(\mathcal{F}) du \leq \frac{c}{n} \int_{\epsilon^2}^{\epsilon} (1/u)^\alpha du \leq \frac{c}{n} (\epsilon^2)^{1-\alpha}.$$

1066 Hence the requirement that this term is bounded by ϵ^2 is equivalent to $\epsilon^{2\alpha} \gtrsim (1/n)$, or $\epsilon^2 \gtrsim (1/n)^{1/\alpha}$.
 1067 When $\alpha \in (1, 2)$, we have $\frac{1}{\alpha} > \frac{2}{2+\alpha}$, so that this condition is milder than our first condition.

1068 Finally, we have $\max \left\{ \frac{\log N_\epsilon(\mathcal{Q})}{n}, \frac{\log N_\epsilon(\Pi)}{n} \right\} \leq \frac{c}{n} (1/\epsilon)^\alpha$, and requiring that this term scales as ϵ^2
 1069 amounts to requiring that $\epsilon^{2+\alpha} \asymp (1/n)$, or equivalently $\epsilon^2 \asymp (1/n)^{\frac{2}{2+\alpha}}$, as before. \square

1070 **C Main Proofs**

1071 This section is devoted to the proofs of our guarantees for general function classes—namely, Proposi-
 1072 tion 4 that holds in a deterministic manner, and Theorem 3 that gives high probability bounds under a
 1073 particular sampling model.

1074 **C.1 Proof of Proposition 4**

1075 Our proof makes use of an elementary simulation lemma, which we state here:

1076 **Lemma 9** (Simulation lemma). *For any policy π and function Q , we have*

$$\mathbb{E}_{S \sim \nu_{\text{start}}}(Q - Q^\pi)(S, \pi) = \frac{\mathbb{E}_\pi \mathcal{B}^\pi Q}{1 - \gamma} \quad (48)$$

1077 See Appendix C.2 for the proof of this claim.

1078 **C.1.1 Proof of policy evaluation claims**

1079 First of all, we have the elementary bounds

$$\begin{aligned} |\widehat{V}_{\min}^\pi - V^\pi| &= \left| \min_{Q \in \widehat{\mathcal{C}}_n^\pi} \mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) - V^\pi \right| \leq \max_{Q \in \widehat{\mathcal{C}}_n^\pi} |\mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) - V^\pi|, \quad \text{and} \\ |\widehat{V}_{\max}^\pi - V^\pi| &= \left| \max_{Q \in \widehat{\mathcal{C}}_n^\pi} \mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) - V^\pi \right| \leq \max_{Q \in \widehat{\mathcal{C}}_n^\pi} |\mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) - V^\pi|. \end{aligned}$$

1080 Consequently, in order to prove the bound (8b) it suffices to upper bound the right-hand side common
 1081 in the two above displays. Since $\widehat{\mathcal{C}}_n^\pi \subseteq \mathcal{C}_n^\pi$, we have the upper bound

$$\begin{aligned} \max_{Q \in \widehat{\mathcal{C}}_n^\pi} |\mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) - V^\pi| &\leq \max_{Q \in \mathcal{C}_n^\pi} |\mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) - V^\pi| \\ &= \max_{Q \in \mathcal{C}_n^\pi} |\mathbb{E}_{S \sim \nu_{\text{start}}} [Q(S, \pi) - Q^\pi(S, \pi)]| \\ &\stackrel{(i)}{=} \frac{1}{1 - \gamma} \max_{Q \in \mathcal{C}_n^\pi} \frac{\mathbb{E}_\pi \mathcal{B}^\pi Q}{1 - \gamma} \end{aligned}$$

1082 where step (i) follows from Lemma 9. Combined with the earlier displays, this completes the proof
 1083 of the bound (8b).

1084 We now show the inclusion $[\widehat{V}_{\min}^\pi, \widehat{V}_{\max}^\pi] \ni V^\pi$ when weak realizability holds. By definition of weak
 1085 realizability, there exists some $Q_\star^\pi \in \mathcal{C}_\infty^\pi$. In conjunction with our sandwich assumption, we are
 1086 guaranteed that $Q_\star^\pi \in \mathcal{C}_\infty^\pi \subseteq \widehat{\mathcal{C}}_n^\pi$, and consequently

$$\begin{aligned} \widehat{V}_{\min}^\pi &= \min_{Q \in \widehat{\mathcal{C}}_n^\pi} \mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) \leq \min_{Q \in \mathcal{C}_\infty^\pi} \mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) \leq \mathbb{E}_{S \sim \nu_{\text{start}}} Q_\star^\pi(S, \pi) = V^\pi, \quad \text{and} \\ \widehat{V}_{\max}^\pi &= \max_{Q \in \widehat{\mathcal{C}}_n^\pi} \mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) \geq \max_{Q \in \mathcal{C}_\infty^\pi} \mathbb{E}_{S \sim \nu_{\text{start}}} Q(S, \pi) \geq \mathbb{E}_{S \sim \nu_{\text{start}}} Q_\star^\pi(S, \pi) = V^\pi. \end{aligned}$$

1087 **C.1.2 Proof of policy optimization claims**

1088 We now prove the oracle inequality (9) on the value $V^{\widetilde{\pi}}$ of a policy $\widetilde{\pi}$ that optimizes the max-min
 1089 criterion. Fix an arbitrary comparator policy π . Starting with the inclusion $[\widehat{V}_{\min}^{\widetilde{\pi}}, \widehat{V}_{\max}^{\widetilde{\pi}}] \ni V^{\widetilde{\pi}}$, we
 1090 have

$$V^{\widetilde{\pi}} \stackrel{(i)}{\geq} \widehat{V}_{\min}^{\widetilde{\pi}} \stackrel{(ii)}{\geq} \widehat{V}_{\min}^\pi = V^\pi - \left(V^\pi - \widehat{V}_{\min}^\pi \right) \stackrel{(iii)}{\geq} V^\pi - \frac{1}{1 - \gamma} \max_{Q \in \mathcal{C}_n^\pi} \frac{|\mathbb{E}_\pi \mathcal{B}^\pi Q|}{1 - \gamma},$$

1091 where step (i) follows from the stated inclusion at the start of the argument; step (ii) follows
 1092 since $\widetilde{\pi}$ solves the max-min program; and step (iii) follows from the bound $|V^\pi - \widehat{V}_{\min}^\pi| \leq$
 1093 $\frac{1}{1 - \gamma} \max_{Q \in \mathcal{C}_n^\pi} \frac{\mathbb{E}_\pi \mathcal{B}^\pi Q}{1 - \gamma}$, as proved in the preceding section. This lower bound holds uniformly
 1094 for all comparators π , from which the stated claim follows.

1095 **C.2 Proof of Lemma 9**

1096 For each $t = 1, 2, \dots$, let \mathbb{E}_t be the expectation over the state-action pair at timestep t upon starting
 1097 from ν_{start} , so that we have $\mathbb{E}_{S \sim \nu_{\text{start}}}(Q - Q^\pi)(S, \pi) = \mathbb{E}_0[Q - Q^\pi]$ by definition. We claim that

$$\mathbb{E}_0[Q - Q^\pi] = \sum_{\tau=1}^t \gamma^{\tau-1} \mathbb{E}_{\tau-1} \mathcal{B}^\pi Q + \gamma^t \mathbb{E}_t[Q - Q^\pi] \quad \text{for all } t = 1, 2, \dots \quad (49)$$

1098 For the base case $t = 1$, we have

$$\mathbb{E}_0[Q - Q^\pi] = \mathbb{E}_0[Q - \mathcal{T}^\pi Q] + \mathbb{E}_0[\mathcal{T}^\pi Q - \mathcal{T}^\pi Q^\pi] = \mathbb{E}_0[Q - \mathcal{T}^\pi Q] + \gamma \mathbb{E}_1[Q - Q^\pi], \quad (50)$$

1099 where we have used the definition of the Bellman evaluation operator to assert that
 1100 $\mathbb{E}_0[\mathcal{T}^\pi Q - \mathcal{T}^\pi Q^\pi] = \gamma \mathbb{E}_1[Q - Q^\pi]$. Since $Q - \mathcal{T}^\pi Q = \mathcal{B}^\pi Q$, the equality (50) is equivalent
 1101 to the claim (49) with $t = 1$.

1102 Turning to the induction step, we now assume that the claim (49) holds for some $t \geq 1$, and show
 1103 that it holds at step $t + 1$. By a similar argument, we can write

$$\begin{aligned} \gamma^t \mathbb{E}_t[Q - Q^\pi] &= \gamma^t \mathbb{E}_t[Q - \mathcal{T}^\pi Q + \mathcal{T}^\pi Q - \mathcal{T}^\pi Q^\pi] = \gamma^t \mathbb{E}_t[Q - \mathcal{T}^\pi Q] + \gamma^{t+1} \mathbb{E}_{t+1}[Q - Q^\pi] \\ &= \gamma^t \mathbb{E}_t \mathcal{B}^\pi Q + \gamma^{t+1} \mathbb{E}_{t+1}[Q - Q^\pi]. \end{aligned}$$

1104 By the induction hypothesis, equality (49) holds for t , and substituting the above equality shows that
 1105 it also holds at time $t + 1$.

1106 Since the equivalence (49) holds for all t , we can take the limit as $t \rightarrow \infty$, and doing so yields the
 1107 claim.

1108 **C.3 Proof of Theorem 3**

1109 In the statement of the theorem, we require choosing $\epsilon > 0$ to satisfy the upper bound $\epsilon^2 \lesssim \frac{\rho(\epsilon, \delta)}{n}$,
 1110 and then provide an upper bound in terms of $\sqrt{\rho(\epsilon, \delta)/n}$. It is equivalent to instead choose ϵ to
 1111 satisfy the lower bound $\epsilon^2 \gtrsim \frac{\rho(\epsilon, \delta)}{n}$, and then provide upper bounds proportional to ϵ . For the
 1112 purposes of the proof, the latter formulation turns out to be more convenient and we pursue it here.
 1113

1114 To streamline notation, let us introduce the shorthand $\langle f, \mathcal{D}^\pi(Q) \rangle \stackrel{\text{def}}{=} \langle f, \delta^\pi(Q) \rangle_n - \langle f, \mathcal{B}^\pi(Q) \rangle_\mu$.
 1115 For each pair (Q, π) , we then define the random variable

$$Z_n(Q, \pi) \stackrel{\text{def}}{=} \sup_{f \in \mathcal{F}^\pi} \frac{|\langle f, \mathcal{D}^\pi(Q) \rangle|}{\sqrt{\|f\|_n^2 + \lambda}}.$$

1116 Central to our proof of the theorem is a uniform bound on this random variable, one that holds for all
 1117 pairs (Q, π) . In particular, our strategy is to exhibit some $\epsilon > 0$ for which, upon setting $\lambda = 4\epsilon^2$, we
 1118 have the guarantees

$$\frac{1}{4} \leq \frac{\sqrt{\|f\|_n^2 + \lambda}}{\sqrt{\|f\|_\mu^2 + \lambda}} \leq 2 \quad \text{uniformly for all } f \in \mathcal{F}, \text{ and} \quad (51a)$$

$$Z_n(Q, \pi) \leq \epsilon \quad \text{uniformly for all } (Q, \pi), \quad (51b)$$

1119 both with probability at least $1 - \delta$. In particular, consistent with the theorem statement, we show
 1120 that this claim holds if we choose $\epsilon > 0$ to satisfy the inequality

$$\epsilon^2 \geq \bar{c} \frac{\rho(\epsilon, \delta)}{n} \quad (52)$$

1121 where $\bar{c} > 0$ is a sufficiently large (but universal) constant.

1122 Supposing that the bounds (51a) and (51b) hold, let us now establish the set inclusions claimed in the
 1123 theorem.

1124 **Inclusion** $\widehat{\mathcal{C}}_\infty^\pi \subseteq \widehat{\mathcal{C}}_n^\pi(\epsilon)$: Define the random variable $M_n(Q, \pi) \stackrel{def}{=} \sup_{f \in \mathcal{F}^\pi} \frac{|\langle f, \mathcal{B}^\pi(Q) \rangle_\mu|}{\sqrt{\|f\|_n^2 + \lambda}}$, and observe
 1125 that $Q \in \widehat{\mathcal{C}}_\infty^\pi$ implies that $M_n(Q, \pi) = 0$. With this definition, we have

$$\sup_{f \in \mathcal{F}^\pi} \frac{|\langle f, \delta^\pi(Q) \rangle_n|}{\sqrt{\|f\|_n^2 + \lambda}} \stackrel{(i)}{\leq} M_n(Q, \pi) + Z_n(Q, \pi) \stackrel{(ii)}{\leq} \epsilon$$

1126 where step (i) follows from the triangle inequality; and step (ii) follows since $M_n(Q, \pi) = 0$, and
 1127 $Z_n(Q, \pi) \leq \epsilon$ from the bound (51b).

1128 **Inclusion** $\widehat{\mathcal{C}}_n^\pi(\epsilon) \subseteq \mathcal{C}_n^\pi(4\epsilon)$ By the definition of $\mathcal{C}_n^\pi(4\epsilon)$, we need to show that

$$\bar{M}(Q, \pi) \stackrel{def}{=} \sup_{f \in \mathcal{F}^\pi} \frac{|\langle f, \mathcal{B}^\pi(Q) \rangle_\mu|}{\sqrt{\|f\|_\mu^2 + \lambda}} \leq 4\epsilon \quad \text{for any } Q \in \widehat{\mathcal{C}}_n^\pi(\epsilon).$$

1129 Now we have

$$\bar{M}(Q, \pi) \stackrel{(i)}{\leq} 2M_n(Q, \pi) \stackrel{(ii)}{\leq} 2 \left\{ \sup_{f \in \mathcal{F}^\pi} \frac{|\langle f, \delta^\pi(Q) \rangle_n|}{\sqrt{\|f\|_n^2 + \lambda}} + Z_n(Q, \pi) \right\} \stackrel{(iii)}{\leq} 2\{\epsilon + \epsilon\} = 4\epsilon,$$

1130 where step (i) follows from the sandwich relation (51a); step (ii) follows from the triangle inequality
 1131 and the definition of $Z_n(Q, \pi)$; and step (iii) follows since $Z_n(Q, \pi) \leq \epsilon$ from the bound (51b), and

$$\sup_{f \in \mathcal{F}^\pi} \frac{|\langle f, \delta^\pi(Q) \rangle_n|}{\sqrt{\|f\|_n^2 + \lambda}} \leq \epsilon, \quad \text{using the inclusion } Q \in \widehat{\mathcal{C}}_n^\pi(\epsilon).$$

1132 Consequently, the remainder of our proof is devoted to establishing the claims (51a) and (51b).
 1133 In doing so, we make repeated use of some Bernstein bounds, stated in terms of the shorthand
 1134 $\Psi_n(\delta) = \frac{\log(n/\delta)}{n}$.

1135 **Lemma 10.** *There is a universal constant c such each the following statements holds with probability
 1136 at least $1 - \delta$. For any f , we have*

$$\left| \|f\|_n^2 - \|f\|_\mu^2 \right| \leq c \left\{ \|f\|_\mu \sqrt{\Psi_n(\delta)} + \Psi_n(\delta) \right\}, \quad (53a)$$

1137 and for any (Q, π) and any function f , we have

$$\left| \langle f, \delta^\pi(Q) \rangle_n - \langle f, \mathcal{B}^\pi(Q) \rangle_\mu \right| \leq c \left\{ \|f\|_\mu \sqrt{\Psi_n(\delta)} + \|f\|_\infty \Psi_n(\delta) \right\}. \quad (53b)$$

1138 These bounds follow by identifying a martingale difference sequence, and applying a form of
 1139 Bernstein's inequality tailored to the martingale setting. See Section C.6.3 for the details.

1140 C.4 Proof of the sandwich relation (51a)

1141 We claim that (modulo the choice of constants) it suffices to show that

$$\left| \|f\|_n - \|f\|_\mu \right| \leq \epsilon \quad \text{uniformly for all } f \in \mathcal{F} \quad (54)$$

1142 for some universal constant c' . Indeed, when this bound holds, we have

$$\|f\|_n + 2\epsilon \leq \|f\|_\mu + 3\epsilon \leq \frac{3}{2} \{ \|f\|_\mu + 2\epsilon \}, \quad \text{and} \quad \|f\|_n + 2\epsilon \geq \|f\|_\mu + \epsilon \geq \frac{1}{2} \{ \|f\|_\mu + 2\epsilon \},$$

1143 so that $\frac{\|f\|_\mu + 2\epsilon}{\|f\|_n + 2\epsilon} \in \left[\frac{1}{2}, \frac{3}{2} \right]$. To relate this statement to the claimed sandwich, observe the inclusion

1144 $\frac{\|f\| + \sqrt{2\epsilon}}{\sqrt{\|f\|^2 + 4\epsilon^2}} \in [1, \sqrt{2}]$, where $\|f\|$ can be either $\|f\|_n$ or $\|f\|_\mu$. Combining this fact with our previous

1145 bound, we see that $\frac{\sqrt{\|f\|_n^2 + 4\epsilon^2}}{\sqrt{\|f\|_\mu^2 + 4\epsilon^2}} \in \left[\frac{1}{\sqrt{2}}, \frac{3\sqrt{2}}{2} \right] \subset \left[\frac{1}{4}, 3 \right]$, as claimed.
 1146

1147 The remainder of our analysis is focused on proving the bound (54). Defining the random variable
 1148 $Y_n(f) = \left| \|f\|_n - \|f\|_\mu \right|$, we need to establish a high probability bound on $\sup_{f \in \mathcal{F}} Y_n(f)$. Let

1149 $\{f^1, \dots, f^N\}$ be an ϵ -cover of \mathcal{F} in the sup-norm. For any $f \in \mathcal{F}$, we can find some f^j such that
 1150 $\|f - f^j\|_\infty \leq \epsilon$, whence

$$\begin{aligned} Y_n(f) &\leq Y_n(f^j) + |Y_n(f^j) - Y_n(f)| \stackrel{(i)}{\leq} Y_n(f^j) + \left| \|f^j\|_n - \|f\|_n \right| + \left| \|f^j\|_\mu - \|f\|_\mu \right| \\ &\stackrel{(ii)}{\leq} Y_n(f^j) + \|f^j - f\|_n + \|f^j - f\|_\mu \\ &\stackrel{(iii)}{\leq} Y_n(f^j) + 2\epsilon, \end{aligned}$$

1151 where steps (i) and (ii) follow from the triangle inequality; and step (iii) follows from the inequality
 1152 $\max\{\|f^j - f\|_n, \|f^j - f\|_\mu\} \leq \|f^j - f\|_\infty \leq \epsilon$. Thus, we have reduced the problem to bounding a
 1153 finite maximum.

1154 Note that if $\max\{\|f^j\|_n, \|f^j\|_\mu\} \leq \epsilon$, then we have $Y_n(f^j) \leq 2\epsilon$ by the triangle inequality. Other-
 1155 wise, we may assume that $\|f^j\|_n + \|f^j\|_\mu \geq \epsilon$. With probability at least $1 - \delta$, we have

$$\begin{aligned} \left| \|f^j\|_n - \|f\|_\mu \right| &= \frac{\left| \|f^j\|_n^2 - \|f\|_\mu^2 \right|}{\|f^j\|_n + \|f\|_\mu} \stackrel{(i)}{\leq} \frac{c\{\|f^j\|_\mu \sqrt{\Psi_n(\delta)} + \Psi_n(\delta)\}}{\|f^j\|_\mu + \|f\|_\mu} \\ &\stackrel{(ii)}{\leq} c\left\{ \sqrt{\Psi_n(\delta)} + \frac{\Psi_n(\delta)}{\epsilon} \right\}, \end{aligned}$$

1156 where step (i) follows from the Bernstein bound (53a) from Lemma 10, and step (ii) uses the fact that
 1157 $\|f^j\|_n + \|f^j\|_\mu \geq \epsilon$.

1158 Taking union bound over all N elements in the cover and replacing δ with δ/N , we have

$$\max_{j \in [N]} Y_n(f^j) \leq c\left\{ \sqrt{\Psi_n(\delta/N)} + \frac{\Psi_n(\delta/N)}{\epsilon} \right\}$$

1159 with probability at least $1 - \delta$. Recalling that $N = N_\epsilon(\mathcal{F})$, our choice (52) of ϵ ensures that
 1160 $\sqrt{\Psi_n(\delta/N)} \leq c\epsilon$ for some universal constant c . Putting together the pieces (and increasing the
 1161 constant \bar{c} in the choice (52) of ϵ as needed) yields the claim.

1162 C.5 Proof of the uniform upper bound (51b)

1163 We need to establish an upper bound on $Z_n(Q, \pi)$ that holds uniformly for all (Q, π) . Our first
 1164 step is to prove a high probability bound for a fixed pair. We then apply a standard discretization
 1165 argument to make it uniform in the pair.

1166 Note that we can write $Z_n(Q, \pi) = \sup_{f \in \mathcal{F}} \frac{V_n(f)}{\sqrt{\|f\|_n^2 + \lambda}}$, where we have defined $V_n(f) \stackrel{def}{=} \langle f, \mathcal{D}^\pi(Q) \rangle$. Our first lemma provides a uniform bound on the latter random variables:

1168 **Lemma 11.** *Suppose that $\epsilon^2 \geq \Psi_n(\delta/N_\epsilon(\mathcal{F}))$. Then we have*

$$V_n(f) \leq c\{\|f\|_\mu \epsilon + \epsilon^2\} \quad \text{for all } f \in \mathcal{F} \quad (55)$$

1169 with probability at least $1 - \delta$.

1170 See Appendix C.6.1 for the proof of this claim.

1171

1172 We claim that the bound (55) implies that, for any fixed pair (Q, π) , we have

$$Y_n(Q, \pi) \leq c'\epsilon \quad \text{with probability at least } 1 - \delta.$$

1173 Indeed, when Lemma 11 holds, for any $f \in \mathcal{F}$, we can write

$$\frac{V_n(f)}{\sqrt{\|f\|_n^2 + \lambda}} = \frac{\sqrt{\|f\|_\mu^2 + \lambda}}{\sqrt{\|f\|_n^2 + \lambda}} \frac{V_n(f)}{\sqrt{\|f\|_\mu^2 + \lambda}} \stackrel{(i)}{\leq} 3 \frac{c\{\|f\|_\mu \epsilon + \epsilon^2\}}{\sqrt{\|f\|_\mu^2 + \lambda}} \stackrel{(ii)}{\leq} c'\epsilon,$$

1174 where step (i) uses the sandwich relation (51a), along with the bound (55); and step (ii) follows given
 1175 the choice $\lambda = 4\epsilon^2$. We have thus proved that for any fixed (Q, π) and $\epsilon \geq \Psi_n(\delta/N_\epsilon(\mathcal{F}))$, we have

$$Z_n(Q, \pi) \leq c'\epsilon \quad \text{with probability at least } 1 - \delta. \quad (56)$$

1176 Our next step is to upgrade this bound to one that is uniform over all pairs (Q, π) . We do so via a
 1177 discretization argument: let $\{Q^j\}_{j=1}^J$ and $\{\pi^k\}_{k=1}^K$ be ϵ -coverings of \mathcal{Q} and Π , respectively.

1178 **Lemma 12.** *We have the upper bound*

$$\sup_{Q, \pi} Z_n(Q, \pi) \leq \max_{(j,k) \in [J] \times [K]} Z_n(Q^j, \pi^k) + 4\epsilon. \quad (57)$$

1179 See Section C.6.2 for the proof of this claim.

1180 If we replace δ with $\delta/(JK)$, then we are guaranteed that the bound (56) holds uniformly over the
 1181 family $\{Q^j\}_{j=1}^J \times \{\pi^k\}_{k=1}^K$. Recalling that $J = N_\epsilon(\mathcal{Q})$ and $K = N_\epsilon(\Pi)$, we conclude that for any
 1182 ϵ satisfying the inequality (52), we have $\sup_{Q, \pi} Z_n(Q, \pi) \leq \tilde{c}\epsilon$ with probability at least $1 - \delta$. (Note
 1183 that by suitably scaling up ϵ via the choice of constant \tilde{c} in the bound (52), we can arrange for $\tilde{c} = 1$,
 1184 as in the stated claim.)

1185 C.6 Proofs of supporting lemmas

1186 In this section, we collect together the proofs of Lemmas 11 and 12, which were stated and used
 1187 in Appendix C.5.

1188 C.6.1 Proof of Lemma 11

1189 We first localize the problem to the class $\mathcal{F}(\epsilon) = \{f \in \mathcal{F} \mid \|f\|_\mu \leq \epsilon\}$. In particular, if there exists
 1190 some $\tilde{f} \in \mathcal{F}$ that violates (55), then the rescaled function $f = \epsilon\tilde{f}/\|\tilde{f}\|_\mu$ belongs to $\mathcal{F}(\epsilon)$, and satisfies
 1191 $V_n(f) \geq c\epsilon^2$. Consequently, it suffices to show that $V_n(f) \leq c\epsilon^2$ for all $f \in \mathcal{F}(\epsilon)$.

1192 Choose an ϵ -cover of \mathcal{F} in the sup-norm with $N = N_\epsilon(\mathcal{F})$ elements. Using this cover, for any
 1193 $f \in \mathcal{F}(\epsilon)$, we can find some f^j such that $\|f - f^j\|_\infty \leq \epsilon$. Thus, for any $f \in \mathcal{F}(\epsilon)$, we can write

$$V_n(f) \leq V_n(f^j) + V_n(f - f^j) \leq \underbrace{V_n(f^j)}_{T_1} + \underbrace{\sup_{g \in \mathcal{G}(\epsilon)} V_n(g)}_{T_2}, \quad (58)$$

1194 where $\mathcal{G}(\epsilon) \stackrel{\text{def}}{=} \{f_1 - f_2 \mid f_1, f_2 \in \mathcal{F}, \|f_1 - f_2\|_\infty \leq \epsilon\}$. We bound each of these two terms in turn.
 1195 In particular, we show that each of T_1 and T_2 are upper bounded by $c\epsilon^2$ with high probability.

1196 **Bounding T_1 :** From the Bernstein bound (53b), we have

$$V_n(f^k) \leq c\{\|f^k\|_\mu \sqrt{\Psi_n(\delta/N)} + \|f^k\|_\infty \Psi_n(\delta/N)\} \quad \text{for all } k \in [N]$$

1197 with probability at least $1 - \delta$. Now for the particular f^j chosen to approximate $f \in \mathcal{F}(\epsilon)$, we have

$$\|f^j\|_\mu \leq \|f^j - f\|_\mu + \|f\|_\mu \leq 2\epsilon,$$

1198 where the inequality follows since $\|f^j - f\|_\mu \leq \|f^j - f\|_\infty \leq \epsilon$, and $\|f\|_\mu \leq \epsilon$. Consequently, we
 1199 conclude that

$$T_1 \leq c\left\{2\epsilon\sqrt{\Psi_n(\delta/N)} + \Psi_n(\delta/N)\right\} \leq c'\epsilon^2 \quad \text{with probability at least } 1 - \delta.$$

1200 where the final inequality follows from our choice of ϵ .

1201 **Bounding T_2 :** Define $\mathcal{G} \stackrel{\text{def}}{=} \{f_1 - f_2 \mid f_1, f_2 \in \mathcal{F}\}$. We need to bound a supremum of the
 1202 process $\{V_n(g), g \in \mathcal{G}\}$ over the subset $\mathcal{G}(\epsilon)$. From the Bernstein bound (53b), the increments
 1203 $V_n(g_1) - V_n(g_2)$ of this process are sub-Gaussian with parameter $\|g_1 - g_2\|_\mu \leq \|g_1 - g_2\|_\infty$, and
 1204 sub-exponential with parameter $\|g_1 - g_2\|_\infty$. Therefore, we can apply a chaining argument that uses
 1205 the metric entropy $\log N_t(\mathcal{G})$ in the supremum norm. Moreover, we can terminate the chaining at 2ϵ ,
 1206 because we are taking the supremum over the subset $\mathcal{G}(\epsilon)$, and it has sup-norm diameter at most 2ϵ .
 1207 Moreover, the lower interval of the chain can terminate at $2\epsilon^2$, since our goal is to prove an upper
 1208 bound of this order. Then, by using high probability bounds for the suprema of empirical processes
 1209 (e.g., Theorem 5.36 in the book [Wai19]), we have

$$T_2 \leq c_1 \int_{2\epsilon^2}^{2\epsilon} \phi\left(\frac{\log N_t(\mathcal{G})}{n}\right) dt + c_2\{\epsilon\sqrt{\Psi_n(\delta)} + \epsilon\Psi_n(\delta)\} + 2\epsilon^2$$

1210 with probability at least $1 - \delta$. (Here the reader should recall our shorthand $\phi(s) = \max\{s, \sqrt{s}\}$.)
 1211 Since \mathcal{G} consists of differences from \mathcal{F} , we have the upper bound $\log N_t(\mathcal{G}) \leq 2 \log N_{t/2}(\mathcal{F})$, and
 1212 hence (after making the change of variable $u = t/2$ in the integrals)

$$T_2 \leq c'_1 \int_{\epsilon^2}^{\epsilon} \phi\left(\frac{\log N_u(\mathcal{F})}{n}\right) du + c_2 \{\epsilon \sqrt{\Psi_n(\delta)} + \epsilon \Psi_n(\delta)\} \leq \tilde{c} \epsilon^2,$$

1213 where the last inequality follows from our choice of ϵ .

1214 C.6.2 Proof of Lemma 12

1215 By our choice of the ϵ -covers, for any (Q, π) , there is a pair (Q^j, π^k) such that

$$\|Q^j - Q\|_{\infty} \leq \epsilon, \quad \text{and} \quad \|\pi^k - \pi\|_{\infty, 1} = \sup_s \|\pi^k(\cdot | s) - \pi(\cdot | s)\|_1 \leq \epsilon.$$

1216 Using this pair, an application of the triangle inequality yields

$$|Z_n(Q, \pi) - Z_n(Q^j, \pi^k)| \leq \underbrace{|Z_n(Q, \pi) - Z_n(Q, \pi^k)|}_{T_1} + \underbrace{|Z_n(Q, \pi^k) - Z_n(Q^j, \pi^k)|}_{T_2}$$

1217 We bound each of these terms in turn, in particular proving that $T_1 + T_2 \leq 24\epsilon$. Putting together the
 1218 pieces yields the bound stated in the lemma.

1219 **Bounding T_2 :** From the definition of Z_n , we have

$$T_2 = |Z_n(Q, \pi^k) - Z_n(Q^j, \pi^k)| \leq \sup_{f \in \mathcal{F}} \frac{|\langle f, \mathcal{D}^{\pi^k}(Q - Q^j) \rangle|}{\sqrt{\|f\|_n^2 + \lambda}}.$$

1220 Now another application of the triangle inequality yields

$$\begin{aligned} |\langle f, \mathcal{D}^{\pi^k}(Q - Q^j) \rangle| &\leq |\langle f, \delta^{\pi^k}(Q - Q^j) \rangle_n| + |\langle f, \mathcal{B}^{\pi^k}(Q - Q^j) \rangle|_{\mu} \\ &\leq \|f\|_n \|\delta^{\pi^k}(Q - Q^j)\|_n + \|f\|_{\mu} \|\mathcal{B}^{\pi^k}(Q - Q^j)\|_{\mu} \\ &\leq \max\{\|f\|_n, \|f\|_{\mu}\} \left\{ \|\delta^{\pi^k}(Q - Q^j)\|_{\infty} + \|\mathcal{B}^{\pi^k}(Q - Q^j)\|_{\infty} \right\} \end{aligned}$$

1221 where step (i) follows from the Cauchy–Schwarz inequality. Now in terms of the shorthand
 1222 $\Delta \stackrel{\text{def}}{=} Q - Q^j$, we have

$$\|\mathcal{B}^{\pi^k}(Q - Q^j)\|_{\infty} = \sup_{(s, a)} \left| \Delta(s, a) - \gamma \mathbb{E}_{s^+ \sim \mathbb{P}(s, a)} [\Delta(s^+, \pi)] \right| \leq 2\|\Delta\|_{\infty} \leq 2\epsilon. \quad (59a)$$

1223 An entirely analogous argument yields

$$\|\delta^{\pi^k}(Q - Q^j)\|_{\infty} \leq 2\epsilon \quad (59b)$$

1224 Conditioned on the sandwich relation (51a), we have $\sup_{f \in \mathcal{F}} \frac{\max\{\|f\|_n, \|f\|_{\mu}\}}{\sqrt{\|f\|_n^2 + \lambda}} \leq 4$. Combining this
 1225 bound with inequalities (59a) and (59b), we have shown that $T_2 \leq 4\{2\epsilon + 2\epsilon\} = 16\epsilon$.

1226 **Bounding T_1 :** In this case, a similar argument yields

$$|\langle f, (\mathcal{D}^{\pi} - \mathcal{D}^{\pi^k})(Q) \rangle| \leq \max\{\|f\|_n, \|f\|_{\mu}\} \left\{ \|(\delta^{\pi} - \delta^{\pi^k})(Q)\|_n + \|(\mathcal{B}^{\pi} - \mathcal{B}^{\pi^k})(Q)\|_{\mu} \right\}.$$

1227 Now we have

$$\begin{aligned} \|(\delta^{\pi} - \delta^{\pi^k})(Q)\|_n &\leq \max_{i=1, \dots, n} \left| \sum_{a'} (\pi(a' | s_i) - \pi^k(a' | s_i)) Q(s_i^+, a') \right| \\ &\leq \max_s \sum_{a'} |\pi(a' | s) - \pi^k(a' | s)| \|Q\|_{\infty} \\ &\leq \epsilon. \end{aligned}$$

1228 A similar argument yields that $\|(\mathcal{B}^{\pi} - \mathcal{B}^{\pi^k})(Q)\|_{\mu} \leq \epsilon$, and arguing as before, we conclude that
 1229 $T_1 \leq 4\{\epsilon + \epsilon\} = 8\epsilon$.

1230 **C.6.3 Proof of Lemma 10**

1231 Our proof of this claim makes use of the following known Bernstein bound for martingale differences
 1232 (cf. Theorem 1 in the paper [BLL⁺11]). Recall the shorthand notation $\Psi_n(\delta) = \frac{\log(n/\delta)}{n}$.

1233 **Lemma 13** (Bernstein's Inequality for Martingales). *Let $\{X_t\}_{t \geq 1}$ be a martingale difference se-*
 1234 *quence with respect to the filtration $\{\mathcal{F}_t\}_{t \geq 1}$. Suppose that $|X_t| \leq 1$ almost surely, and let \mathbb{E}_t denote*
 1235 *expectation conditional on \mathcal{F}_t . Then for all $\delta \in (0, 1)$, we have*

$$\left| \frac{1}{n} \sum_{t=1}^n X_t \right| \leq 2 \left[\left(\frac{1}{n} \sum_{t=1}^n \mathbb{E}_t X_t^2 \right) \Psi_n(2\delta) \right]^{1/2} + 2\Psi_n(2\delta) \quad (60)$$

1236 with probability at least $1 - \delta$.

1237 With this result in place, we divide our proof into two parts, corresponding to the two claims (53b)
 1238 and (53a) stated in Lemma 10.

1239 **Proof of the bound (53b):** Recall that at step i , the triple (s, a, o) is drawn according to a condi-
 1240 tional distribution $\mu_i(\cdot \mid \mathcal{F}_i)$. Similarly, we let d_i denote the distribution of (s, a, r, s^+, o) conditioned
 1241 on the filtration \mathcal{F}_i . Note that μ_i is obtained from d_i by marginalizing out the pair (r, s^+) . Moreover,
 1242 by the tower property of expectation, the Bellman error is equivalent to the average TD error.

1243 Using these facts, we have the equivalence

$$\begin{aligned} \langle f, \delta^\pi Q \rangle_{d_i} &= \mathbb{E}_{d_i} \{ f(s, a, o) [Q(s, a) - r - \gamma Q(s^+, \pi)] \} \\ &= \mathbb{E}_{(s, a, o) \sim \mu_i} \{ f(s, a, o) \mathbb{E}_{r \sim R(s, a), s^+ \sim \mathbb{P}(s, a)} [Q(s, a) - r - \gamma Q(s^+, \pi)] \} \\ &= \mathbb{E}_{(s, a, o) \sim \mu_i} \{ f(s, a, o) [Q(s, a) - (\mathcal{T}^\pi Q)(s, a)] \} \\ &= \langle f, \mathcal{B}^\pi Q \rangle_{\mu_i}. \end{aligned}$$

1244 As a consequence, we can write $\langle f, \delta^\pi(Q) \rangle_n - \langle f, \mathcal{B}^\pi(Q) \rangle_\mu = \frac{1}{n} \sum_{i=1}^n W_i$ where

$$W_i \stackrel{\text{def}}{=} f(s_i, a_i, o_i) [Q(s_i, a_i) - r_i - \gamma Q(s_i^+, \pi)] - \mathbb{E}_{d_i} \{ f(s, a, o) [Q(s, a) - r - \gamma Q(s^+, \pi)] \}$$

1245 defines a martingale difference sequence (MDS). Thus, we can prove the claim by applying a
 1246 Bernstein martingale inequality.

1247 Since $\|r\|_\infty \leq 1$ and $\|Q\|_\infty \leq 1$ by assumption, we have $\|W_i\|_\infty \leq 3\|f\|_\infty$, and

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{d_i} [W_i^2] \leq 9 \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mu_i} [f^2(s_i, a_i, o_i)] = 9\|f\|_\mu^2.$$

1248 Consequently, the claimed bound (53b) follows by applying the Bernstein bound stated in Lemma 13.

1249 **Proof of the bound (53a):** In this case, we have the additive decomposition

$$\|f\|_n^2 - \|f\|_\mu^2 = \frac{1}{n} \sum_{i=1}^n \underbrace{\left\{ f^2(s_i, a_i, o_i) - \mathbb{E}_{\mu_i} [f^2(s, a, o)] \right\}}_{W'_i},$$

1250 where $\{W'_i\}_{i=1}^n$ again defines a martingale difference sequence. Note that $\|W'_i\|_\infty \leq 2\|f\|_\infty^2 \leq 2$,
 1251 and

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mu_i} [(W'_i)^2] \stackrel{(i)}{\leq} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mu_i} [f^4(S, A, O)] \leq \|f\|_\infty^2 \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mu_i} [f^2(S, A, O)] \stackrel{(ii)}{\leq} \|f\|_\mu^2,$$

1252 where step (i) uses the fact that the variance of f^2 is at most the fourth moment, and step (ii) uses the
 1253 bound $\|f\|_\infty \leq 1$. Consequently, the claimed bound (53a) follows by applying the Bernstein bound
 1254 stated in Lemma 13.

1255 **D Proofs for Section 4 and Appendix A.5**

1256 In this section, we collect together the proofs of results stated without proof in Section 4 and
1257 Appendix A.5.

1258 **D.1 Proof of Proposition 1**

1259 *Proof.* Since $f^* \in \mathcal{F}^\pi$, we are guaranteed that the corresponding constraint must hold. It reads as

$$|\mathbb{E}_\mu \frac{1}{b_\pi} \frac{d_\pi}{\mu} \mathcal{B}^\pi Q|^2 = \frac{1}{b_\pi^2} |\mathbb{E}_\pi \mathcal{B}^\pi Q|^2 \stackrel{(iii)}{\leq} \left(\frac{1}{b_\pi^2} \|\frac{d_\pi}{\mu}\|_\mu^2 + \lambda \right) \frac{\rho}{n}.$$

1260 where step (iii) follows from the definition of population constraint. Re-arranging yields the upper
1261 bound

$$\frac{|\mathbb{E}_\mu \frac{d_\pi}{\mu} \mathcal{B}^\pi Q|^2}{(1 + \lambda) \frac{\rho}{n}} \leq \frac{(\|\frac{d_\pi}{\mu}\|_\mu^2 + b_\pi^2 \lambda) \frac{\rho}{n}}{(1 + \lambda) \frac{\rho}{n}} = \frac{\mathbb{E}_\pi \left[\frac{d_\pi(S, A)}{\mu(S, A)} \right] + b_\pi^2 \lambda}{1 + \lambda},$$

1262 where the final step uses the fact that

$$\|\frac{d_\pi}{\mu}\|_\mu^2 = \mathbb{E}_\mu \frac{d_\pi^2(S, A)}{\mu^2(S, A)} = \mathbb{E}_\pi \frac{d_\pi(S, A)}{\mu(S, A)}$$

1263 Thus, we have established the bound (i) in our claim (12).

1264 The upper bound (ii) follows immediately since $\mathbb{E}_\pi \frac{d_\pi(s, a)}{\mu(s, a)} \leq \sup_{(s, a)} \frac{d_\pi(s, a)}{\mu(s, a)} \leq b_\pi$.

1265 □

1266 **D.2 Proof of Lemma 1**

1267 Some simple algebra yields

$$\mathcal{B}^\pi Q - \mathcal{B}^\pi Q_*^\pi = [Q - \mathcal{T}^\pi Q] - [Q_*^\pi - \mathcal{T}^\pi Q_*^\pi] = (\mathcal{I} - \gamma \mathbb{P}^\pi)(Q - Q_*^\pi) = (\mathcal{I} - \gamma \mathbb{P}^\pi)\epsilon.$$

1268 Taking expectations under π and recalling that $\langle f, \mathcal{B}^\pi Q_*^\pi \rangle_\pi = 0$ for all $f \in \mathcal{F}^\pi$ yields

$$\langle f, \mathcal{B}^\pi Q \rangle_\pi = \langle f, (\mathcal{I} - \gamma \mathbb{P}^\pi)\epsilon \rangle_\pi.$$

1269 Notice that for any $Q \in \mathcal{Q}^\pi$ there exists a test function $\epsilon = Q - Q_*^\pi \in \mathcal{E}^\pi$, and the associated
1270 population constraint reads

$$\frac{|\langle \epsilon, (\mathcal{I} - \gamma \mathbb{P}^\pi)\epsilon \rangle_\mu|}{\sqrt{\|\epsilon\|_\mu^2 + \lambda}} \leq \sqrt{\frac{\rho}{n}}.$$

1271 Consequently, the off-policy cost coefficient can be upper bounded as

$$K^\pi \leq \max_{\epsilon \in \mathcal{E}_*^\pi} \left\{ \frac{\rho}{n} \frac{\langle \mathbf{1}, (\mathcal{I} - \gamma \mathbb{P}^\pi)\epsilon \rangle_\pi^2}{1 + \lambda} \right\} \leq \max_{\epsilon \in \mathcal{E}_*^\pi} \left\{ \frac{\|\epsilon\|_\mu^2 + \lambda}{\|\mathbf{1}\|_\pi^2 + \lambda} \frac{\langle \mathbf{1}, (\mathcal{I} - \gamma \mathbb{P}^\pi)\epsilon \rangle_\pi^2}{\langle \epsilon, (\mathcal{I} - \gamma \mathbb{P}^\pi)\epsilon \rangle_\mu^2} \right\},$$

1272 as claimed in the bound (14).

1273 **D.3 Proof of Lemma 7**

1274 If weak Bellman closure holds, then we can write

$$\mathcal{B}^\pi Q = Q - \mathcal{T}^\pi Q = Q - \mathcal{P}^\pi(Q) \in \mathcal{E}^\pi.$$

1275 For any $Q \in \mathcal{Q}^\pi$, the function $\epsilon = Q - \mathcal{P}^\pi(Q)$ belongs to \mathcal{E}^π , and the associated population constraint
1276 reads $\frac{|\langle \epsilon, \epsilon \rangle_\mu|}{\sqrt{\|\epsilon\|_\mu^2 + \lambda}} \leq \sqrt{\frac{\rho}{n}}$. Consequently, the off-policy cost coefficient is upper bounded as

$$K^\pi \leq \max_{\epsilon \in \mathcal{E}^\pi} \left\{ \frac{n}{\rho} \frac{v \langle \mathbf{1}, \epsilon \rangle_\pi^2}{1 + \lambda} \right\} \leq \max_{\epsilon \in \mathcal{E}^\pi} \left\{ \frac{\|\epsilon\|_\mu^2 + \lambda}{1 + \lambda} \frac{\langle \mathbf{1}, \epsilon \rangle_\pi^2}{\langle \epsilon, \epsilon \rangle_\mu^2} \right\} \leq \max_{\epsilon \in \mathcal{E}^\pi} \left\{ \frac{\langle \mathbf{1}, \epsilon \rangle_\pi^2}{\langle \epsilon, \epsilon \rangle_\mu^2} \right\},$$

1277 where the final inequality follows from the fact that $\|\epsilon\|_\mu \leq 1$.

1278 **D.4 Proof of Lemma 2**

1279 We split our proof into the two separate claims.

1280 **Proof of the bound (16a):** When the test function class includes $\mathcal{F}_\pi^{\mathcal{B}}$, then any Q feasible must
 1281 satisfy the population constraints

$$\frac{\langle \mathcal{B}^\pi Q', \mathcal{B}^\pi Q \rangle_\mu}{\sqrt{\|\mathcal{B}^\pi Q'\|_\mu^2 + \lambda}} \leq \sqrt{\frac{\rho}{n}}, \quad \text{for all } Q' \in \mathcal{Q}^\pi.$$

1282 Setting $Q' = Q$ yields $\frac{\|\mathcal{B}^\pi Q\|_\mu^2}{\sqrt{\|\mathcal{B}^\pi Q\|_\mu^2 + \lambda}} \leq \sqrt{\frac{\rho}{n}}$. If $\|\mathcal{B}^\pi Q\|_\mu^2 \geq \lambda$, then the claim holds, given our choice

1283 $\lambda = c \frac{\rho}{n}$ for some constant c . Otherwise, the constraint can be weakened to $\frac{\|\mathcal{B}^\pi Q\|_\mu^2}{\sqrt{2\|\mathcal{B}^\pi Q\|_\mu^2}} \leq \sqrt{\frac{\rho}{n}}$, which
 1284 yields the bound (16a).

1285 **Proof of the bound (16b):** We now prove the sequence of inequalities stated in equation (16b).
 1286 Inequality (i) follows directly from the definition of K^π and Lemma 2. Turning to inequality (ii), an
 1287 application of Jensen's inequality yields

$$\langle \mathbf{1}, \mathcal{B}^\pi Q \rangle_\pi^2 = [\mathbb{E}_\pi \mathcal{B}^\pi Q]^2 \leq \mathbb{E}_\pi [\mathcal{B}^\pi Q]^2 = \|\mathcal{B}^\pi Q\|_\pi^2.$$

1288 Finally, inequality (iii) follows by observing that

$$\sup_{Q \in \mathcal{Q}^\pi} \frac{\|\mathcal{B}^\pi Q\|_\pi^2}{\|\mathcal{B}^\pi Q\|_\mu^2} = \sup_{Q \in \mathcal{Q}^\pi} \frac{\mathbb{E}_\pi [(\mathcal{B}^\pi Q)(s, a)]^2}{\mathbb{E}_\mu [(\mathcal{B}^\pi Q)(s, a)]^2} = \sup_{Q \in \mathcal{Q}^\pi} \frac{\mathbb{E}_\mu \left[\frac{d_\pi(s, a)}{\mu(s, a)} \right] [(\mathcal{B}^\pi Q)(s, a)]^2}{\mathbb{E}_\mu [(\mathcal{B}^\pi Q)(s, a)]^2} \leq \sup_{(s, a)} \frac{d_\pi(s, a)}{\mu(s, a)}.$$

1289 **E Proofs for the Linear Setting**

1290 We now prove the results stated in Section 5. Throughout this section, the reader should recall that Q
 1291 takes the linear function $Q(s, a) = \langle w, \phi(s, a) \rangle$, so that the bulk of our arguments operate directly
 1292 on the weight vector $w \in \mathbb{R}^d$.

1293 Given the linear structure, the population and empirical covariance matrices of the feature vectors
 1294 play a central role. We make use of the following known result (cf. Lemma 1 in the paper [ZJZ21])
 1295 that relates these objects:

1296 **Lemma 14** (Covariance Concentration). *There are universal constants (c_1, c_2, c_3) such that for any*
 1297 $\delta \in (0, 1)$, *we have*

$$c_1 \mathbb{E}_\mu \phi \phi^\top \preceq \frac{1}{n} \sum_{i=1}^n \phi_i \phi_i^\top + \frac{c_2}{n} \log \frac{nd}{\delta} I \preceq c_3 \mathbb{E}_\mu \phi \phi^\top + \frac{c_4}{n} \log \frac{nd}{\delta} I. \quad (61)$$

1298 *with probability at least $1 - \delta$.*

1299 **E.1 Proof of Proposition 2**

1300 Under weak realizability, we have

$$\langle f_j, \mathcal{B}^\pi Q_\star^\pi \rangle_\mu = 0 \quad \text{for all } j = 1, \dots, d. \quad (62)$$

1301 Thus, at (s, a) the Bellman error difference reads

$$\begin{aligned} \mathcal{B}^\pi Q(s, a) - \mathcal{B}^\pi Q_\star^\pi(s, a) &= [Q - \mathcal{T}^\pi Q](s, a) - [Q_\star^\pi - \mathcal{T}^\pi Q_\star^\pi](s, a) \\ &= [Q - Q_\star^\pi](s, a) - \gamma \mathbb{E}_{s^+ \sim \mathbb{P}(s, a)} [Q - Q_\star^\pi](s^+, \pi) \\ &= \langle w - w_\star^\pi, \phi(s, a) - \gamma \phi^{+\pi}(s, a) \rangle \end{aligned} \quad (63)$$

1302 To proceed we need the following auxiliary result:

1303 **Lemma 15** (Linear Parameter Constraints). *With probability at least $1 - \delta$, there exists a universal*
 1304 *constant $c_1 > 0$ such that if $Q \in \mathcal{C}_n^\pi$ then $\|w - w_\star^\pi\|_{\Sigma_{\lambda, \text{Boot}}^{+\pi}}^2 \leq c_1 \frac{d\rho}{n}$.*

1305 See Appendix E.2 for the proof.

1306 Using this lemma, we can bound the OPC coefficient as follows

$$\begin{aligned} K^\pi &\stackrel{(i)}{\leq} \frac{n}{\rho} \max_{Q \in \mathcal{C}_n^\pi} \langle \mathbb{1}, \mathcal{B}^\pi Q - \mathcal{B}^\pi Q_\star^\pi \rangle_\pi^2 \stackrel{(ii)}{\leq} \frac{n}{\rho} [\mathbb{E}_\pi (\phi - \gamma \phi^{+\pi})^\top (w - w_\star^\pi)]^2 \\ &\leq \frac{n}{\rho} \|\mathbb{E}_\pi \phi - \gamma \phi^{+\pi}\|_{(\Sigma_{\lambda, \text{Boot}}^{+\pi})^{-1}}^2 \|w - w_\star^\pi\|_{\Sigma_{\lambda, \text{Boot}}^{+\pi}}^2 \\ &\leq c_1 d \|\mathbb{E}_\pi \phi - \gamma \phi^{+\pi}\|_{(\Sigma_{\lambda, \text{Boot}}^{+\pi})^{-1}}^2. \end{aligned}$$

1307 Here step (i) follows from the definition of off-policy cost coefficient, (ii) leverages the linear
 1308 structure and (iii) is Cauchy-Schwartz.

1309 **E.2 Proof of Lemma 15**

1310 Under the event of Theorem 3, the statement of Eq. (51a) holds, and in particular

$$\frac{1}{c_1(\sqrt{\|f\|_\mu^2 + \lambda})} \geq \frac{1}{\sqrt{\|f\|_n^2 + \lambda}} \geq \frac{1}{c_2(\sqrt{\|f\|_\mu^2 + \lambda})}.$$

1311 Thus, the j constraint reads

$$\frac{L}{\sqrt{n}} \gtrsim \frac{\langle f_j, \mathcal{B}^\pi Q \rangle_\mu}{\sqrt{\|f\|_n^2 + \lambda}} = \frac{\langle f_j, \mathcal{B}^\pi Q \rangle_\mu}{\sqrt{\hat{\lambda}_j + \lambda}}$$

1312 where the last step follows from

$$\|f_j\|_{\mathcal{D}}^2 = \frac{1}{n} \sum_{(s,a,r,s^+) \in \mathcal{D}} (f_j(s,a))^2 = \frac{1}{n} \sum_{i=1}^n (\hat{u}_j^\top \phi_i)^2 = \hat{u}_j^\top \widehat{\Sigma} \hat{u}_j = \hat{\lambda}_j.$$

1313 Now, squaring and summing over the constraints and using Eq. (63) yields

$$\begin{aligned} d \frac{L^2}{n} &\gtrsim \sum_{j=1}^m \left\langle \frac{\hat{u}_j^\top \phi}{\sqrt{\hat{\lambda}_j + \lambda}}, (\phi - \gamma \phi^{+\pi})^\top (w - w_\star^\pi) \right\rangle_\mu^2 \\ &= \sum_{j=1}^m \left[\frac{\hat{u}_j^\top}{\sqrt{\hat{\lambda}_j + \lambda}} \mathbb{E}_\mu \phi (\phi - \gamma \phi^{+\pi})^\top (w - w_\star^\pi) \right]^2 \\ &= \sum_{j=1}^m \left[\frac{\hat{u}_j^\top}{\sqrt{\hat{\lambda}_j + \lambda}} \underbrace{(\Sigma - \gamma \Sigma^{+\pi})(w - w_\star^\pi)}_{\stackrel{\text{def}}{=} y} \right]^2 \\ &= y^\top \left(\sum_{j=1}^m \frac{\hat{u}_j \hat{u}_j^\top}{\hat{\lambda}_j + \lambda} \right) y \\ &= y^\top \left(\widehat{\Sigma} + \lambda I \right)^{-1} y \\ &\gtrsim y^\top \Sigma_\lambda^{-1} y. \end{aligned}$$

1314 The last inequality holds via Lemma 14 (*Covariance Concentration*) with probability at least $1 - \delta$
1315 since λ is a large enough regularizer. Let us complete the quadratic form:

$$\|y + \lambda(w - w_\star^\pi)\|_{\Sigma_\lambda^{-1}}^2 \leq (\|y\|_{\Sigma_\lambda^{-1}} + \lambda\|(w - w_\star^\pi)\|_{\Sigma_\lambda^{-1}})^2 \lesssim \|y\|_{\Sigma_\lambda^{-1}}^2 + \lambda.$$

1316 Therefore, adding λ to both sides of the prior display and noticing that $\lambda \lesssim \frac{L^2}{n}$ gives

$$\begin{aligned} d \frac{L^2}{n} &\gtrsim \|y + \lambda(w - w_\star^\pi)\|_{\Sigma_\lambda^{-1}}^2 \\ &= (w - w_\star^\pi)(\Sigma_\lambda - \gamma \Sigma^{+\pi})^\top (\Sigma_\lambda^{-1}) (\Sigma_\lambda - \gamma \Sigma^{+\pi})(w - w_\star^\pi) \\ &= (w - w_\star^\pi)(\Sigma_{\lambda, \text{Boot}}^{+\pi})(w - w_\star^\pi) \\ &= \|(w - w_\star^\pi)\|_{\Sigma_{\lambda, \text{Boot}}^{+\pi}}^2. \end{aligned}$$

1317 E.3 Proof of Proposition 3

1318 Under weak Bellman closure, we have

$$\mathcal{B}^\pi Q = Q - \mathcal{T}^\pi Q = \phi^\top (w - \mathcal{P}^\pi(w)). \quad (64)$$

1319 With a slight abuse of notation, let $\mathcal{P}^\pi(w)$ denote the weight vector that defines the action-value
1320 function $\mathcal{P}^\pi(Q)$. We introduce the following auxiliary lemma:

1321 **Lemma 16** (Linear Parameter Constraints with Bellman Closure). *With probability at least $1 - \delta$, if*
1322 *$Q \in \mathcal{C}_n^\pi$ then $\|w - \mathcal{P}^\pi(w)\|_{\Sigma_\lambda}^2 \leq c_1 \frac{d\rho}{n}$.*

1323 See Appendix E.4 for the proof. Using this lemma, we can bound the OPC coefficient as follows

$$\begin{aligned} K^\pi &\stackrel{(i)}{\leq} \frac{n}{\rho} \max_{Q \in \mathcal{C}_n^\pi} \langle \mathbb{1}, \mathcal{B}^\pi Q \rangle_\pi^2 \stackrel{(ii)}{\leq} \frac{n}{\rho} [\mathbb{E}_\pi(\phi)^\top (w - \mathcal{P}^\pi(w))]^2 \\ &\stackrel{(iii)}{\leq} \frac{n}{\rho} \|\mathbb{E}_\pi \phi\|_{(\Sigma_\lambda)^{-1}}^2 \|w - \mathcal{P}^\pi(w)\|_{\Sigma_\lambda}^2 \\ &\leq c_1 d \|\mathbb{E}_\pi \phi\|_{(\Sigma_\lambda)^{-1}}^2. \end{aligned}$$

1324 Here step (i) follows from the definition of off-policy cost coefficient, (ii) leverages the linear
1325 structure and (iii) is Cauchy-Schwartz.

1326 **E.4 Proof of Appendix E.4**

1327 Under the event of Theorem 3, the statement of Eq. (51a) holds, and in particular

$$\frac{1}{c_1(\sqrt{\|f\|_\mu^2 + \lambda})} \geq \frac{1}{\sqrt{\|f\|_n^2 + \lambda}} \geq \frac{1}{c_2(\sqrt{\|f\|_\mu^2 + \lambda})}.$$

1328 Thus, the j constraint reads

$$\frac{L}{\sqrt{n}} \gtrsim \frac{\langle f_j, \mathcal{B}^\pi Q \rangle_\mu}{\sqrt{\|f\|_n^2 + \lambda}} = \frac{\langle f_j, \mathcal{B}^\pi Q \rangle_\mu}{\sqrt{\hat{\lambda}_j + \lambda}}$$

1329 where the last step follows from

$$\|f_j\|_{\mathcal{D}}^2 = \frac{1}{n} \sum_{(s,a,r,s^+) \in \mathcal{D}} (f_j(s,a))^2 = \frac{1}{n} \sum_{i=1}^n (\hat{u}_j^\top \phi_i)^2 = \hat{u}_j^\top \hat{\Sigma} \hat{u}_j = \hat{\lambda}_j.$$

1330 Now, squaring and summing over the constraints and using Eq. (64) yields

$$\begin{aligned} d \frac{L^2}{n} &\gtrsim \sum_{j=1}^m \left\langle \frac{\hat{u}_j^\top \phi}{\sqrt{\hat{\lambda}_j + \lambda}}, \phi^\top (w - \mathcal{P}^\pi(w)) \right\rangle_\mu^2 \\ &= \sum_{j=1}^m \left[\frac{\hat{u}_j^\top}{\sqrt{\hat{\lambda}_j + \lambda}} \mathbb{E}_\mu \phi \phi^\top (w - \mathcal{P}^\pi(w)) \right]^2 \\ &= \sum_{j=1}^m \left[\frac{\hat{u}_j^\top}{\sqrt{\hat{\lambda}_j + \lambda}} \underbrace{\Sigma(w - \mathcal{P}^\pi(w))}_{\stackrel{\text{def}}{=} y} \right]^2 \\ &= y^\top \left(\sum_{j=1}^m \frac{\hat{u}_j \hat{u}_j^\top}{\hat{\lambda}_j + \lambda} \right) y \\ &= y^\top (\hat{\Sigma} + \lambda I)^{-1} y \\ &\gtrsim y^\top \Sigma_\lambda^{-1} y. \end{aligned}$$

1331 The last inequality holds via Lemma 14 (*Covariance Concentration*) with probability at least $1 - \delta$
 1332 since λ is a large enough regularizer. Let us complete the quadratic form:

$$\|y + \lambda(w - \mathcal{P}^\pi(w))\|_{\Sigma_\lambda^{-1}}^2 \leq (\|y\|_{\Sigma_\lambda^{-1}} + \lambda \|w - \mathcal{P}^\pi(w)\|_{\Sigma_\lambda^{-1}})^2 \lesssim \|y\|_{\Sigma_\lambda^{-1}}^2 + \lambda.$$

1333 Therefore, adding λ to both sides of the prior display and noticing that $\lambda \lesssim \frac{L^2}{n}$ gives

$$\begin{aligned} d \frac{L^2}{n} &\gtrsim \|y + \lambda(w - \mathcal{P}^\pi(w))\|_{\Sigma_\lambda^{-1}}^2 \\ &= (w - \mathcal{P}^\pi(w)) \Sigma_\lambda^\top \left(\Sigma_\lambda^{-1} \right) \Sigma_\lambda (w - \mathcal{P}^\pi(w)) \\ &= (w - \mathcal{P}^\pi(w)) (\Sigma_\lambda) (w - \mathcal{P}^\pi(w)) \\ &= \|(w - \mathcal{P}^\pi(w))\|_{\Sigma_\lambda}^2. \end{aligned}$$

1334 **F Proof of Theorem 2**

1335 In this section, we prove the guarantee on our actor-critic procedure stated in Theorem 2.

1336 **F.1 Adversarial MDPs**

1337 We now introduce sequence of adversarial MDPs $\{\mathcal{M}_t\}_{t=1}^T$ used in the analysis. Each MDP \mathcal{M}_t
 1338 is defined by the same state-action space and transition law as the original MDP \mathcal{M} , but with the
 1339 reward functions R perturbed by R_t —that is

$$\mathcal{M}_t \stackrel{def}{=} \langle \mathcal{S}, \mathcal{A}, R + R_t, \mathbb{P}, \gamma \rangle. \quad (65)$$

1340 For an arbitrary policy π , we denote with Q_t^π and with A_t^π the action value function and the advantage
 1341 function on \mathcal{M}_t ; the value of π from the starting distribution ν_{start} is denoted by V_t^π . We immediately
 1342 have the following expression for the value function, which follows because the dynamics of \mathcal{M}_t and
 1343 \mathcal{M} are identical and the reward function of \mathcal{M}_t equals that of \mathcal{M} plus R_t

$$V_t^\pi \stackrel{def}{=} \frac{1}{1-\gamma} \mathbb{E}_\pi [R + R_t]. \quad (66)$$

1344 Consider the action value function \widehat{Q}_{π_t} returned by the critic, and let the reward perturbation
 1345 $R_t = \mathcal{B}^{\pi_t} \widehat{Q}_{\pi_t}$ be the Bellman error of the critic value function \widehat{Q}_{π_t} . The special property of \mathcal{M}_t is
 1346 that the action value function of π_t on \mathcal{M}_t equals the critic lower estimate \widehat{Q}_{π_t} .

1347 **Lemma 17** (Adversarial MDP Equivalence). *Given the perturbed MDP \mathcal{M}_t from equation (65) with*
 1348 *$R_t \stackrel{def}{=} \mathcal{B}^{\pi_t} \widehat{Q}_{\pi_t}$, we have the equivalence*

$$Q_t^{\pi_t} = \widehat{Q}_{\pi_t}.$$

1349 *Proof.* We need to check that \widehat{Q}_{π_t} solves the Bellman evaluation equations for the adversarial MDP,
 1350 ensuring that \widehat{Q}_{π_t} is the action-value function of π_t on \mathcal{M}_t . Let $\mathcal{T}_t^{\pi_t}$ be the Bellman evaluation
 1351 operator on \mathcal{M}_t for policy π_t . We have

$$\widehat{Q}_{\pi_t} - \mathcal{T}_t^{\pi_t}(\widehat{Q}_{\pi_t}) = \widehat{Q}_{\pi_t} - \mathcal{T}^{\pi_t}(\widehat{Q}_{\pi_t}) - R_t = \mathcal{B}^{\pi_t} \widehat{Q}_{\pi_t} - \mathcal{B}^{\pi_t} \widehat{Q}_{\pi_t} = 0.$$

1352 Thus, the function \widehat{Q}_{π_t} is the action value function of π_t on \mathcal{M}_t , and it is by definition denoted by
 1353 $Q_t^{\pi_t}$. \square

1354 This lemma shows that the action-value function \widehat{Q}_{π_t} computed by the critic is equivalent to the
 1355 action-value function of π_t on \mathcal{M}_t . Thus, we can interpret the critic as performing a model-based
 1356 pessimistic estimate of π_t ; this view is useful in the rest of the analysis.

1357 **F.2 Equivalence of Updates**

1358 The second step is to establish the equivalence between the update rule (22), or equivalently as the
 1359 update (67a), to the exponentiated gradient update rule (67b).

1360 **Lemma 18** (Equivalence of Updates). *For linear Q -functions of the form $Q_t(s, a) = \langle w_t, \phi(s, a) \rangle$,*
 1361 *the parameter update*

$$\pi_{t+1}(a | s) \propto \exp(\phi(s, a)^\top (\theta_t + \eta w_t)), \quad (67a)$$

1362 *is equivalent to the policy update*

$$\pi_{t+1}(a | s) \propto \pi_t(a | s) \exp(\eta Q_t(s, a)), \quad \pi_1(a | s) = \frac{1}{|\mathcal{A}_s|}. \quad (67b)$$

1363 *Proof.* We prove this claim via induction on t . The base case ($t = 1$) holds by a direct calculation.
 1364 Now let us show that the two update rules update π_t in the same way. As an inductive step, assume
 1365 that both rules maintain the same policy $\pi_t \propto \exp(\phi(s, a)^\top \theta_t)$ at iteration t ; we will show the
 1366 policies are still the same at iteration $t + 1$. At any (s, a) , we have

$$\begin{aligned} \pi_{t+1}(a | s) &\propto \exp(\phi(s, a)^\top (\theta_t + \eta w_t)) \propto \exp(\phi(s, a)^\top \theta_t) \exp(\eta \phi(s, a)^\top w_t) \\ &\propto \pi_t(a | s) \exp(\eta Q_t(s, a)). \end{aligned}$$

1367 □

1368 Recall that θ_t is the parameter associated to π_t and that w_t is the parameter associated to \widehat{Q}_{π_t} . Using
 1369 Lemma 18 together with Lemma 17 we obtain that the actor policy π_t satisfies through its parameter
 1370 θ_t the mirror descent update rule (67b) with $Q_t = \widehat{Q}_{\pi_t} = Q_t^{\pi_t}$ and $\pi_1(a | s) = 1/|\mathcal{A}_s|$, $\forall (s, a)$. In
 1371 words, the actor is using Mirror descent to find the best policy on the sequence of adversarial MDPs
 1372 $\{\mathcal{M}_t\}$ implicitly identified by the critic.

1373 F.3 Mirror Descent on Adversarial MDPs

1374 Our third step is to analyze the behavior of mirror descent on the MDP sequence $\{\mathcal{M}_t\}_{t=1}^T$, and then
 1375 translate such guarantees back to the original MDP \mathcal{M} . The following result provides a bound on the
 1376 average of the value functions $\{V^{\pi_t}\}_{t=1}^T$ induced by the actor's policy sequence. This bound involves
 1377 a form of optimization error⁸ given by

$$\mathcal{E}_{opt}(T) = 2 \sqrt{\frac{2 \log |\mathcal{A}|}{T}},$$

1378 as is standard in mirror descent schemes. It also involves the *perturbed rewards* given by $R_t \stackrel{def}{=} \mathcal{B}^{\pi_t} Q_t^{\pi_t}$.

1380 **Lemma 19** (Mirror Descent on Adversarial MDPs). *For any positive integer T , applying the update*
 1381 *rule (67b) with $Q_t = Q_t^{\pi_t}$ for T rounds yields a sequence such that*

$$\frac{1}{T} \sum_{t=1}^T [V^{\tilde{\pi}} - V^{\pi_t}] \leq \frac{1}{1-\gamma} \left\{ \mathcal{E}_{opt}(T) + \frac{1}{T} \sum_{t=1}^T [-\mathbb{E}_{\tilde{\pi}} R_t + \mathbb{E}_{\pi_t} R_t] \right\}, \quad (68)$$

1382 *valid for any comparator policy $\tilde{\pi}$.*

1383 See Appendix F.6 for the proof.

1384

1385 To be clear, the comparator policy $\tilde{\pi}$ need belong to the soft-max policy class. Apart from the
 1386 optimization error term, our bound (68) involves the behavior of the perturbed rewards R_t along the
 1387 comparator $\tilde{\pi}$ and π_t , respectively. These correction terms arise because the actor performs the policy
 1388 update using the action-value function $Q_t^{\pi_t}$ on the perturbed MDPs instead of the real underlying
 1389 MDP.

1390 F.4 Pessimism: Bound on $\mathbb{E}_{\pi_t} R_t$

1391 The fourth step of the proof is to leverage the pessimistic estimates returned by critic to simplify
 1392 equation (68). Using Lemma 9 and the definition of adversarial reward R_t we can write

$$\widehat{V}_{\min}^{\pi} - V^{\pi_t} = \frac{1}{1-\gamma} \langle \mathbb{1}, \mathcal{B}^{\pi_t} \widehat{Q}_{\pi_t} \rangle_{\pi_t} = \frac{1}{1-\gamma} \mathbb{E}_{\pi_t} \mathcal{B}^{\pi_t} \widehat{Q}_{\pi_t} = \frac{1}{1-\gamma} \mathbb{E}_{\pi_t} R_t.$$

1393 Since weak realizability holds, Theorem 3 guarantees that $\widehat{V}_{\min}^{\pi} \leq V^{\pi}$ uniformly for all $\pi \in \Pi$ with
 1394 probability at least $1 - \delta$. Coupled with the prior display, we find that

$$\mathbb{E}_{\pi_t} R_t \leq 0. \quad (69)$$

1395 Using the above display, the result in Eq. (68) can be further upper bounded and simplified.

⁸Technically, this error should depend on $|\mathcal{A}_s|$, if we were to allow the action spaces to have varyign cardinality, but we elide this distinction here.

1396 **F.5 Concentrability: Bound on $\mathbb{E}_{\tilde{\pi}} R_t$**

1397 The term $\mathbb{E}_{\tilde{\pi}} R_t$ can be interpreted as an approximate concentrability factor for the approximate
1398 algorithm that we are investigating.

1399 **Bound under only weak realizability:** Lemma 15 gives with probability at least $1 - \delta$ that any
1400 surviving Q in $\mathcal{C}_n^{\pi_t}$ must satisfy: $\|w - w_{\star}^{\pi_t}\|_{\Sigma_{\lambda, \text{Boot}}^{+\pi_t}}^2 \lesssim \frac{d\rho}{n}$ where $w_{\star}^{\pi_t}$ is the parameter associated to
1401 the weak solution $Q_{\star}^{\pi_t}$. Such bound must apply to the parameter $w_t \in \widehat{\mathcal{C}}_n^{\pi_t}$ identified by the critic.⁹
1402 We are now ready to bound the remaining adversarial reward along the distribution of the comparator
1403 $\tilde{\pi}$.

$$\begin{aligned} |\mathbb{E}_{\tilde{\pi}} R_t| &= |\mathbb{E}_{\tilde{\pi}} \mathcal{B}^{\pi_t} \widehat{Q}_{\pi_t}| \\ &\stackrel{(i)}{=} |\mathbb{E}_{\tilde{\pi}} (\phi - \gamma \phi^{+\pi_t})^\top (w_t - w_{\star}^{\pi_t})| \\ &\leq \|\mathbb{E}_{\tilde{\pi}} [\phi - \gamma \phi^{+\pi_t}]\|_{(\Sigma_{\lambda, \text{Boot}}^{+\pi_t})^{-1}} \|w_t - w_{\star}^{\pi_t}\|_{\Sigma_{\lambda, \text{Boot}}^{+\pi_t}} \\ &\leq c \sqrt{\frac{d\rho}{n}} \sup_{\pi \in \Pi} \left\{ \|\mathbb{E}_{\tilde{\pi}} [\phi - \gamma \phi^{+\pi}]\|_{(\Sigma_{\lambda, \text{Boot}}^{+\pi})^{-1}} \right\}. \end{aligned} \quad (70)$$

1404 Step (i) follows from the expression (63) for the weak Bellman error, along with the definition of the
1405 weak solution $Q_{\star}^{\pi_t}$.

1406 **Bound under weak Bellman closure:** When Bellman closure holds we proceed analogously. The
1407 bound in Lemma 16 ensures with probability at least $1 - \delta$ that $\|w - \mathcal{P}^{\pi_t}(w)\|_{\Sigma_{\lambda}}^2 \leq c \frac{d\rho}{n}$ for all
1408 $w \in \mathcal{C}_n^{\pi_t}$; as before, this relation must apply to the parameter chosen by the critic $w_t \in \widehat{\mathcal{C}}_n^{\pi_t}$. The
1409 bound on the adversarial reward along the distribution of the comparator $\tilde{\pi}$ now reads

$$\begin{aligned} |\mathbb{E}_{\tilde{\pi}} R_t| &= |\mathbb{E}_{\tilde{\pi}} \mathcal{B}^{\pi_t} \widehat{Q}_{\pi_t}| \stackrel{(i)}{=} |\mathbb{E}_{\tilde{\pi}} \phi^\top (w_t - \mathcal{P}^{\pi_t}(w_t))| \\ &\leq \|\mathbb{E}_{\tilde{\pi}} \phi\|_{\Sigma_{\lambda}^{-1}} \|w_t - \mathcal{P}^{\pi_t}(w_t)\|_{\Sigma_{\lambda}} \\ &\leq c \|\mathbb{E}_{\tilde{\pi}} \phi\|_{\Sigma_{\lambda}^{-1}} \sqrt{\frac{d\rho}{n}}. \end{aligned} \quad (71)$$

1410 Here step (i) follows from the expression (64) for the Bellman error under weak closure.

1411 **F.6 Proof of Lemma 19**

1412 We now prove our guarantee for a mirror descent procedure on the sequence of adversarial MDPs.
1413 Our analysis makes use of a standard result on online mirror descent for linear functions (e.g., see
1414 Section 5.4.2 of Hazan [Haz21]), which we state here for reference. Given a finite cardinality set
1415 \mathcal{X} , a function $f : \mathcal{X} \rightarrow \mathbb{R}$, and a distribution ν over \mathcal{X} , we define $f(\nu) \stackrel{\text{def}}{=} \sum_{x \in \mathcal{X}} \nu(x) f(x)$. The
1416 following result gives a guarantee that holds uniformly for any sequence of functions $\{f_t\}_{t=1}^T$, thereby
1417 allowing for the possibility of adversarial behavior.

1418 **Proposition 5** (Adversarial Guarantees for Mirror Descent). *Suppose that we initialize with the*
1419 *uniform distribution $\nu_1(x) = \frac{1}{|\mathcal{X}|}$ for all $x \in \mathcal{X}$, and then perform T rounds of the update*

$$\nu_{t+1}(x) \propto \nu_t(x) \exp(\eta f_t(x)), \quad \text{for all } x \in \mathcal{X}, \quad (72)$$

1420 *using $\eta = \sqrt{\frac{\log |\mathcal{X}|}{2T}}$. If $\|f_t\|_{\infty} \leq 1$ for all $t \in [T]$ then we have the bound*

$$\frac{1}{T} \sum_{t=1}^T [f_t(\tilde{\nu}) - f_t(\nu_t)] \leq \mathcal{E}_{\text{opt}}(T) \stackrel{\text{def}}{=} 2\sqrt{\frac{2 \log |\mathcal{X}|}{T}}. \quad (73)$$

1421 *where $\tilde{\nu}$ is any comparator distribution over \mathcal{X} .*

⁹We abuse the notation and write $w \in \widehat{\mathcal{C}}_n^{\pi}$ in place of $Q \in \widehat{\mathcal{C}}_n^{\pi}$

1422 We now use this result to prove our claim. So as to streamline the presentation, it is convenient to
 1423 introduce the advantage function corresponding to π_t . It is a function of the state-action pair (s, a)
 1424 given by

$$A_t^{\pi_t}(s, a) \stackrel{def}{=} Q_t^{\pi_t}(s, a) - \mathbb{E}_{a^+ \sim \pi_t(\cdot | s)} Q_t^{\pi_t}(s, a^+).$$

1425 In the sequel, we omit dependence on (s, a) when referring to this function, consistent with the rest
 1426 of the paper.

1427 From our earlier observation (66), recall that the reward function of the perturbed MDP \mathcal{M}_t corre-
 1428 sponds to that of \mathcal{M} plus the perturbation R_t . Combining this fact with a standard simulation lemma
 1429 (e.g., [K⁺03]) applied to \mathcal{M}_t , we find that

$$V^{\tilde{\pi}} - V^{\pi_t} = V_t^{\tilde{\pi}} - V_t^{\pi_t} + \frac{1}{1-\gamma} \left[-\mathbb{E}_{\tilde{\pi}} R_t + \mathbb{E}_{\pi_t} R_t \right] = \frac{1}{1-\gamma} \left[\mathbb{E}_{\tilde{\pi}} A_t^{\pi_t} - \mathbb{E}_{\tilde{\pi}} R_t + \mathbb{E}_{\pi_t} R_t \right]. \quad (74a)$$

1430 Now for any given state s , we introduce the linear objective function

$$f_t(\nu) \stackrel{def}{=} \mathbb{E}_{a \sim \nu} Q_t^{\pi_t}(s, a) = \sum_{a \in \mathcal{A}} \nu(a) Q_t^{\pi_t}(s, a),$$

1431 where ν is a distribution over the action space. With this choice, we have the equivalence

$$\mathbb{E}_{a \sim \tilde{\pi}} A_t^{\pi_t}(s, a) = f_t(\tilde{\pi}(\cdot | s)) - f_t(\pi_t(\cdot | s)),$$

1432 where the reader should recall that we have fixed an arbitrary state s . Consequently, applying the
 1433 bound (73) with $\mathcal{X} = \mathcal{A}$ and these choices of linear functions, we conclude that

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{a \sim \tilde{\pi}} A_t^{\pi_t}(s, a) \leq \mathcal{E}_{opt}(T). \quad (74b)$$

1434 This bound holds for any state, and also for any average over the states.

1435 We now combine the pieces to conclude. By computing the average of the bound (74a) over all T
 1436 iterations, we find that

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \left[V^{\tilde{\pi}} - V^{\pi_t} \right] &\leq \frac{1}{1-\gamma} \left\{ \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\tilde{\pi}} A_t^{\pi_t} + \frac{1}{T} \sum_{t=1}^T \left[-\mathbb{E}_{\tilde{\pi}} R_t + \mathbb{E}_{\pi_t} R_t \right] \right\} \\ &\leq \frac{1}{1-\gamma} \left\{ \mathcal{E}_{opt}(T) + \frac{1}{T} \sum_{t=1}^T \left[-\mathbb{E}_{\tilde{\pi}} R_t + \mathbb{E}_{\pi_t} R_t \right] \right\}, \end{aligned}$$

1437 where the final inequality follow from the bound (73), applied for each s . We have thus established
 1438 the claim.